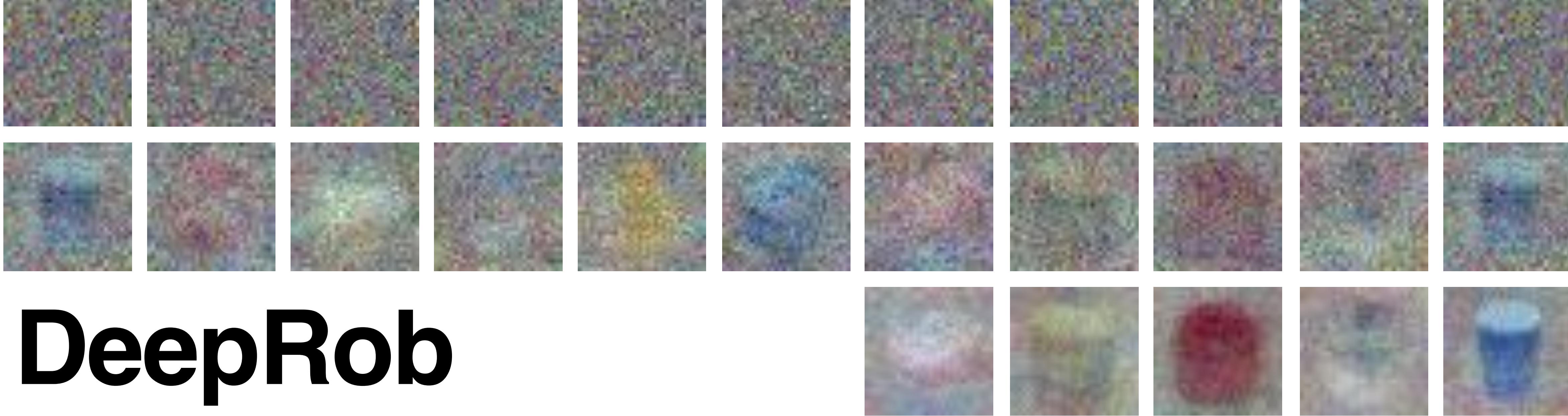


**DR**



# DeepRob

**Lecture 4**  
**Regularization + Optimization**  
**University of Michigan and University of Minnesota**



# Project 1 – Reminder

---

- Instructions and code available on the website
  - Here: [deeprob.org/projects/project1/](https://deeprob.org/projects/project1/)
- Uses Python, PyTorch and Google Colab
- Implement KNN, linear SVM, and linear softmax classifiers
- **Autograder is online**
- **Due Thursday, January 26th 11:59 PM EST**



# Project 1 – Dataset

## Progress Robot Object Perception Samples Dataset



**10 classes**  
**32x32 RGB images**  
**50k training images (5k per class)**  
**10k test images (1k per class)**

Chen et al., “ProgressLabeller: Visual Data Stream Annotation for Training Object-Centric 3D Perception”, IROS, 2022.

# Discussion 2—How was this dataset created?

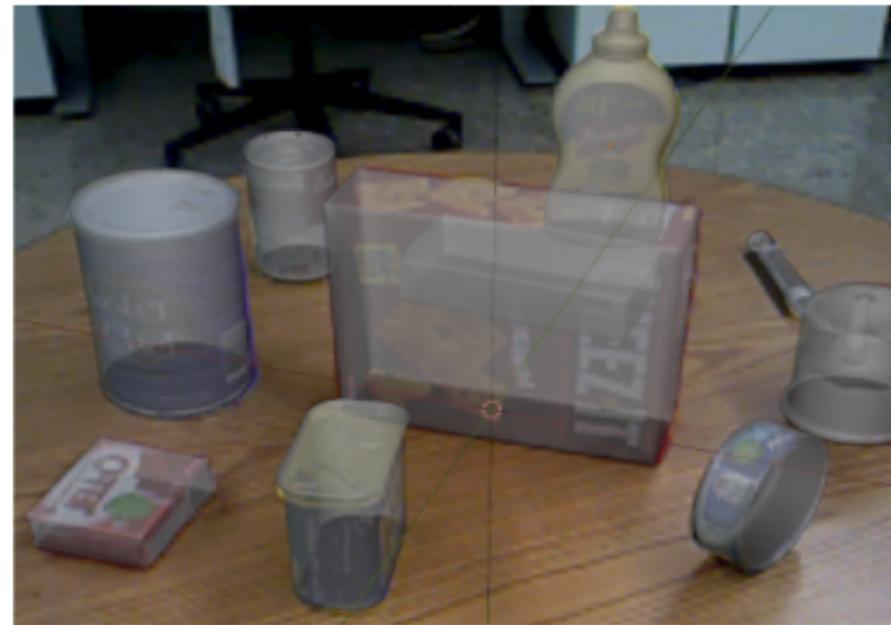
**ProgressLabeller: Visual Data Stream Annotation for Training  
Object-Centric 3D Perception**

Xiaotong Chen Huijie Zhang Zeren Yu Stanley Lewis Odest Chadwicke Jenkins

**Rough Pose Estimates  
from Pretrained Model**



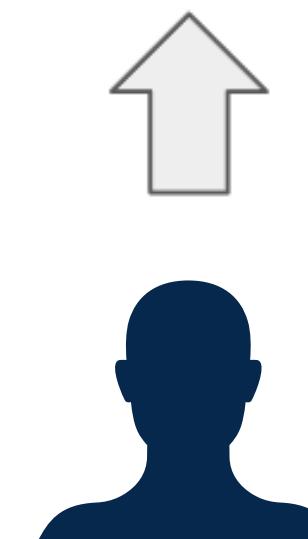
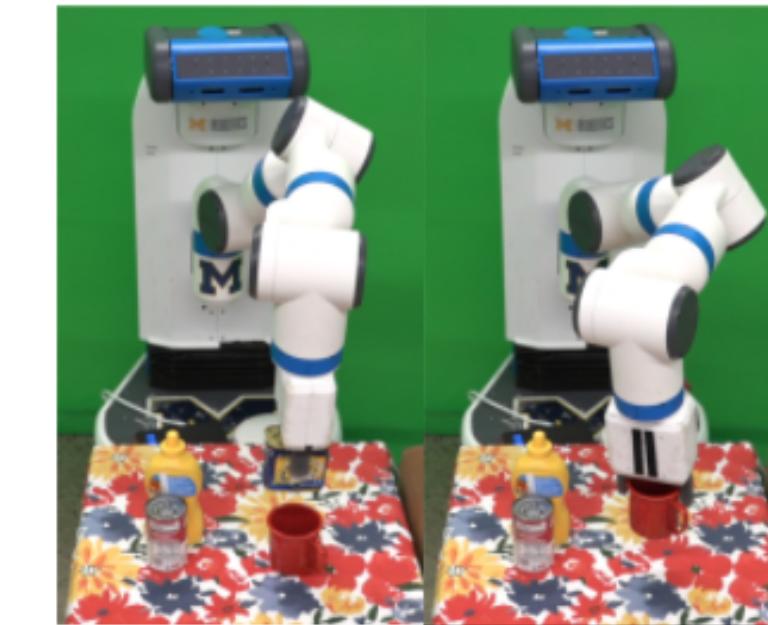
**6D pose annotation through  
interactive interface**



**Fine-tuned Pose  
Estimates**



**Pose-based Robot  
Grasping**



**Human  
Annotator**

**Idea:**

1. Record video of scene
2. Human labels object pose in selected frames
3. Pose labels propagate to (large number of) remaining frames

# Gradescope Quizzes

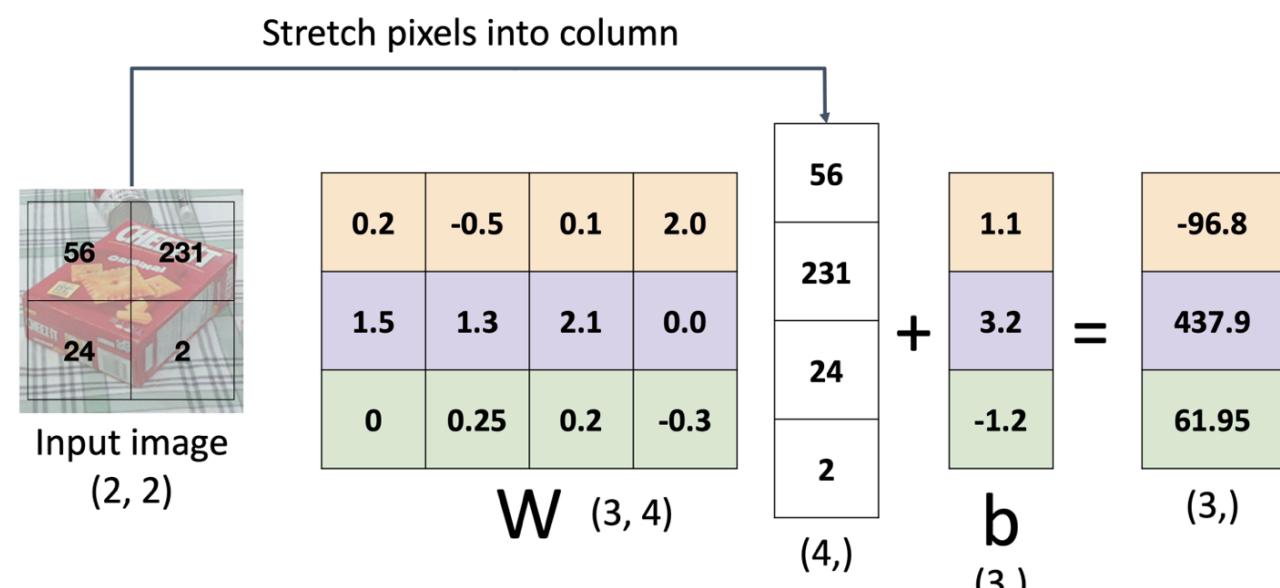
---

- Let course staff know if you have issues accessing
- Quiz links available through gradescope course 480760
- Time limit of 15 min once quiz is opened
- Each available to take from 7:00AM–3:00PM EST on quiz days
- Covers material from previous lectures and graded projects
- **Today only: quiz 1 available until 6:00PM EST**

# Recap—Linear Classifiers

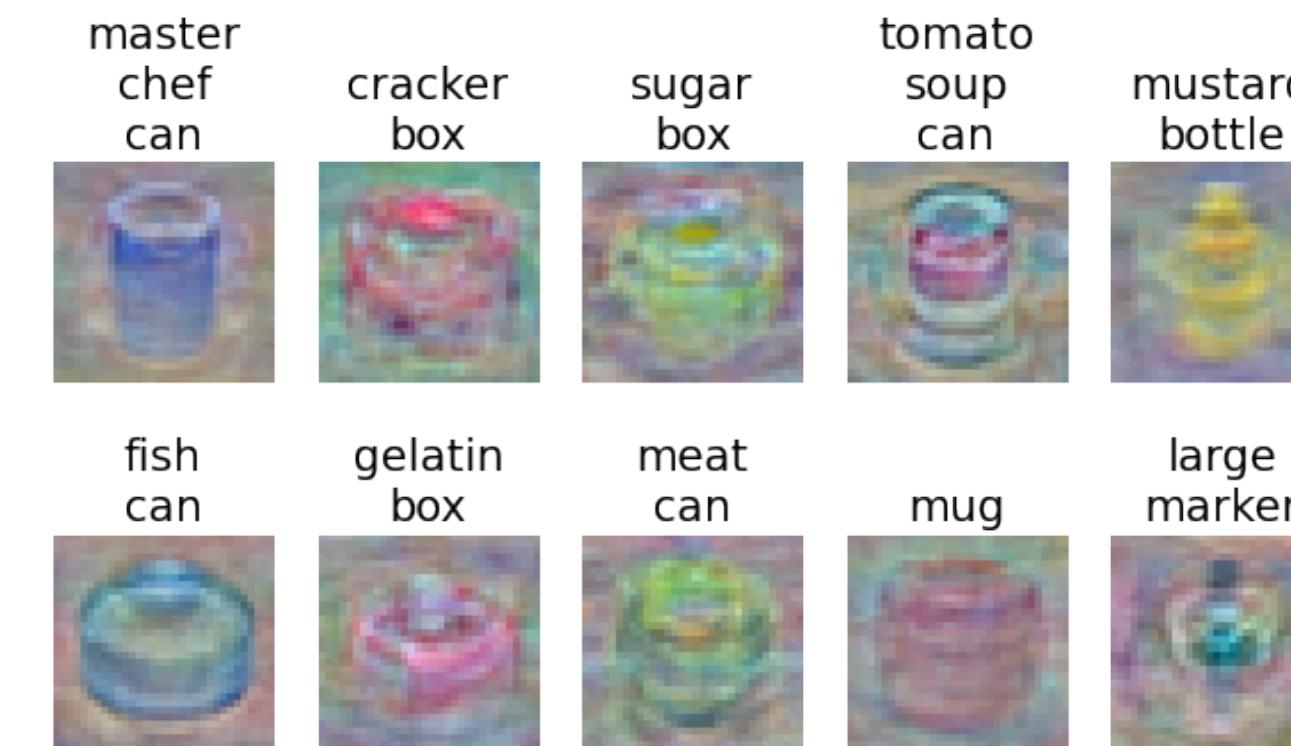
## Algebraic Viewpoint

$$f(x, W) = Wx$$



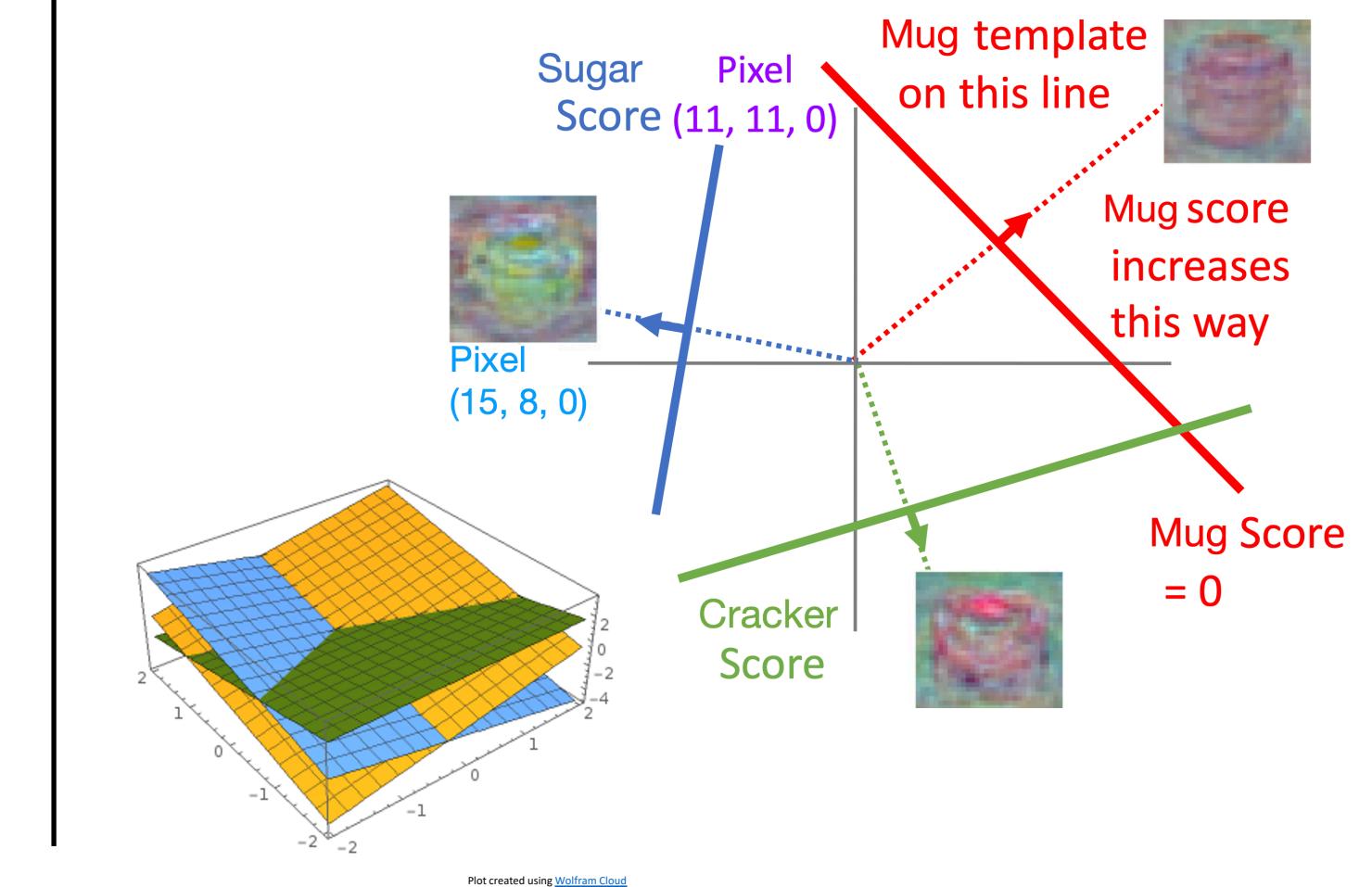
## Visual Viewpoint

One template per class



## Geometric Viewpoint

Hyperplanes cutting up space



# Recap—Loss Functions Quantify Preferences

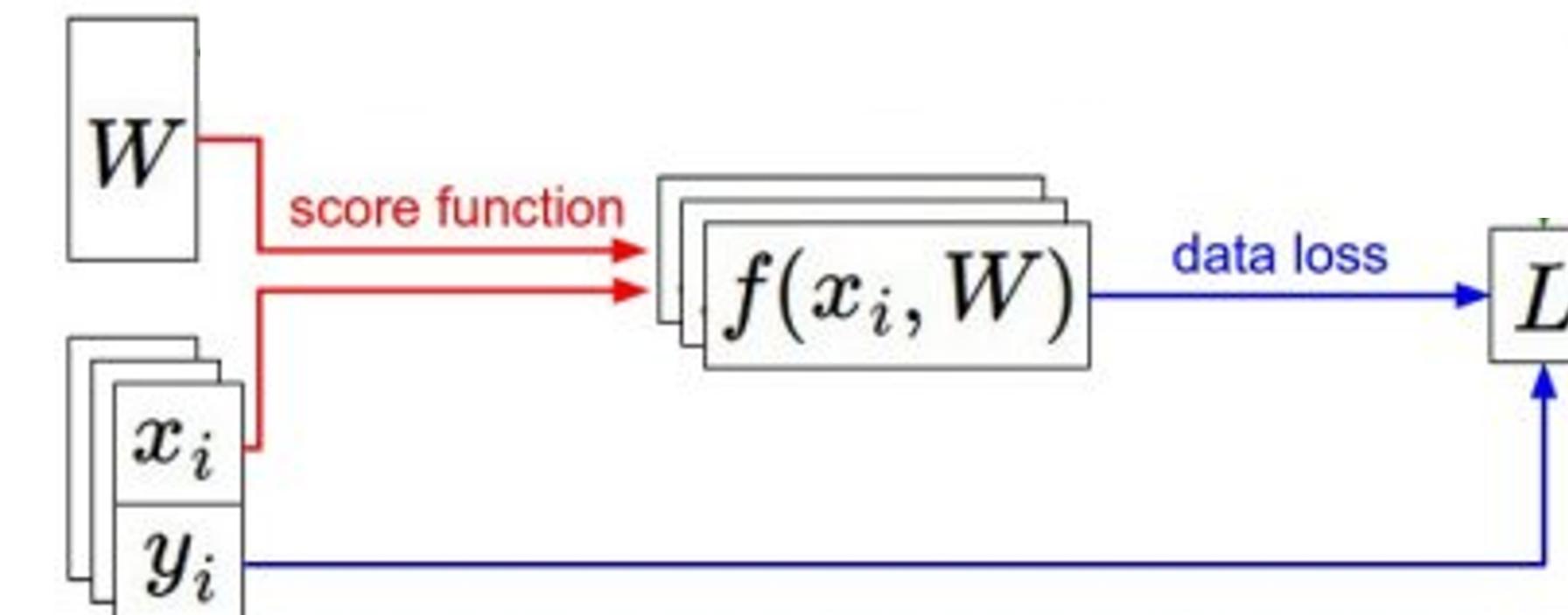
- We have some dataset of  $(x, y)$
- We have a **score function**:
- We have a **loss function**:

$$s = f(x; W, b) = Wx + b$$

Linear classifier

**Softmax:**  $L_i = -\log \left( \frac{\exp(s_{y_i})}{\sum_j \exp(s_j)} \right)$

**SVM:**  $L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$



# Recap—Loss Functions Quantify Preferences

---

**Q: How do we find the best  $W, b$ ?**

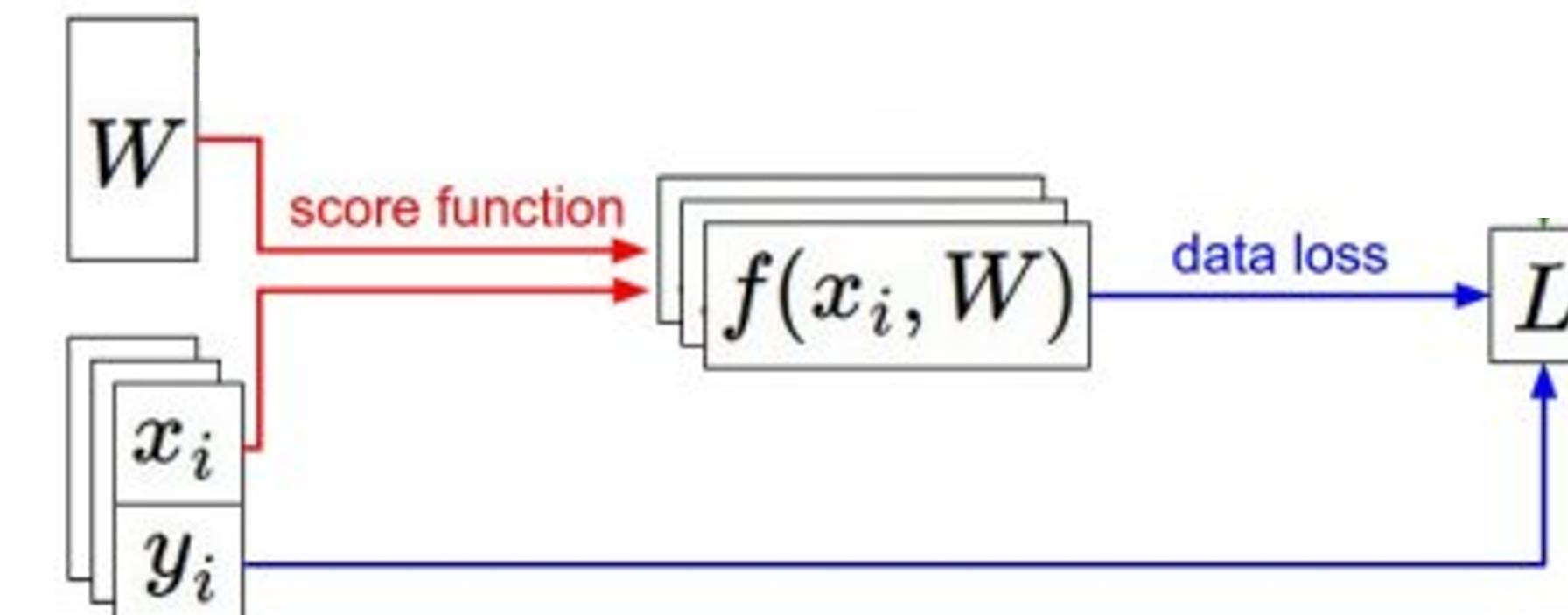
- We have some dataset of  $(x, y)$
- We have a **score function**:
- We have a **loss function**:

$$s = f(x; W, b) = Wx + b$$

Linear classifier

**Softmax:**  $L_i = -\log \left( \frac{\exp(s_{y_i})}{\sum_j \exp(s_j)} \right)$

**SVM:**  $L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$



# Recap—Loss Functions Quantify Preferences

---

- We have some dataset of  $(x, y)$
- We have a **score function**:
- We have a **loss function**:

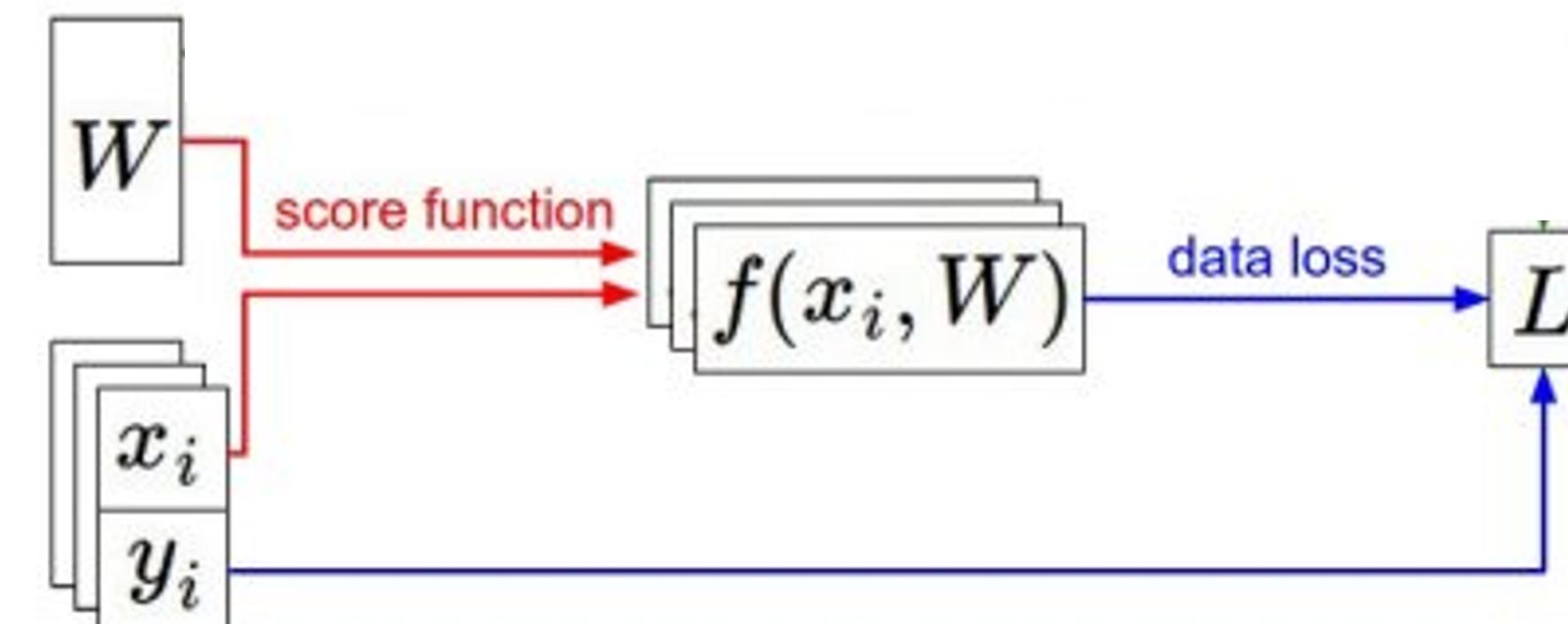
**Problem:** Loss functions encourage good performance on training data but we care about test data

$$s = f(x; W, b) = Wx + b$$

Linear classifier

**Softmax:**  $L_i = -\log \left( \frac{\exp(s_{y_i})}{\sum_j \exp(s_j)} \right)$

**SVM:**  $L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$



# Regularization + Optimization

# Overfitting

---

A model is overfit when it performs too well on the training data, and has poor performance for unseen data

**Both models have perfect accuracy on the training data!**



# Overfitting

---

A model is overfit when it performs too well on the training data, and has poor performance for unseen data

Example: Linear classifier with 1D inputs, 2 classes, and softmax loss

$$s_i = w_i x + b_i$$

$$p_i = \frac{\exp(s_i)}{\exp(s_1) + \exp(s_2)}$$

$$L = -\log(p_y)$$

**Both models have perfect accuracy on the training data!**



# Overfitting

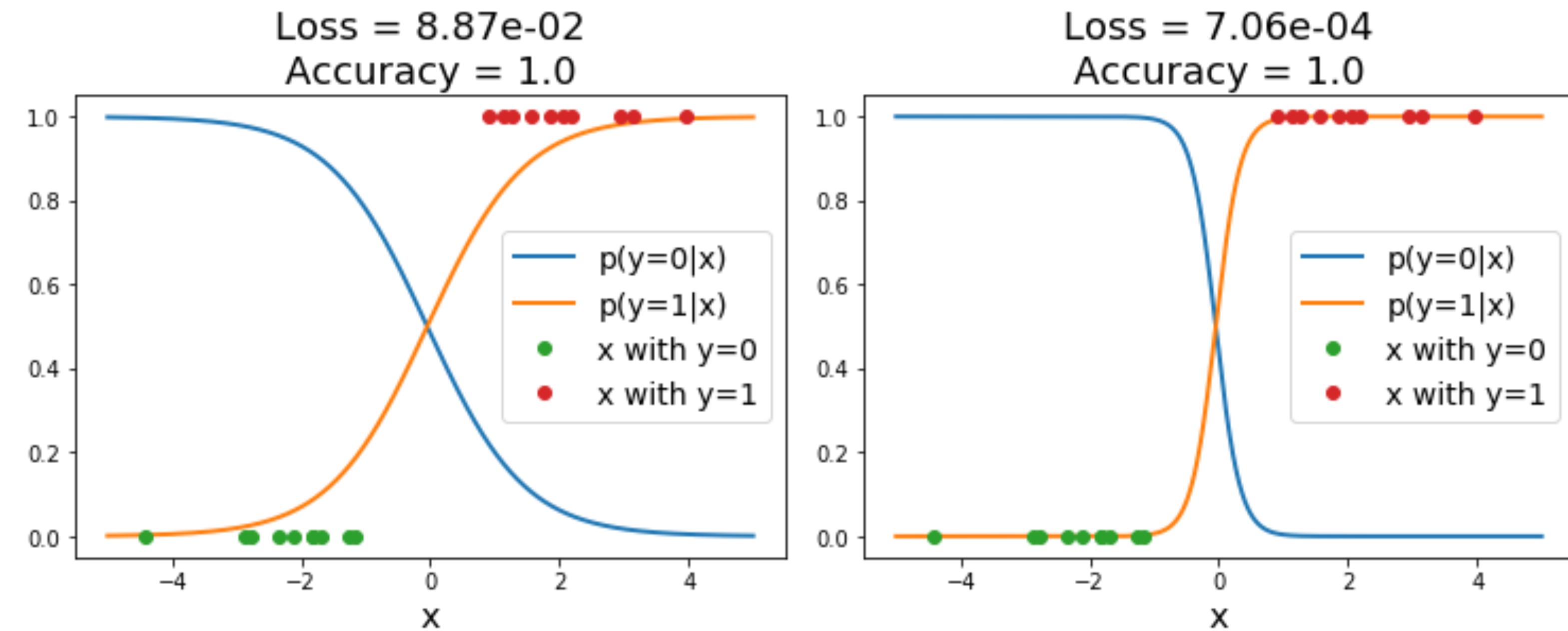
A model is overfit when it performs too well on the training data, and has poor performance for unseen data

Example: Linear classifier with 1D inputs, 2 classes, and softmax loss

$$s_i = w_i x + b_i$$

$$p_i = \frac{\exp(s_i)}{\exp(s_1) + \exp(s_2)}$$

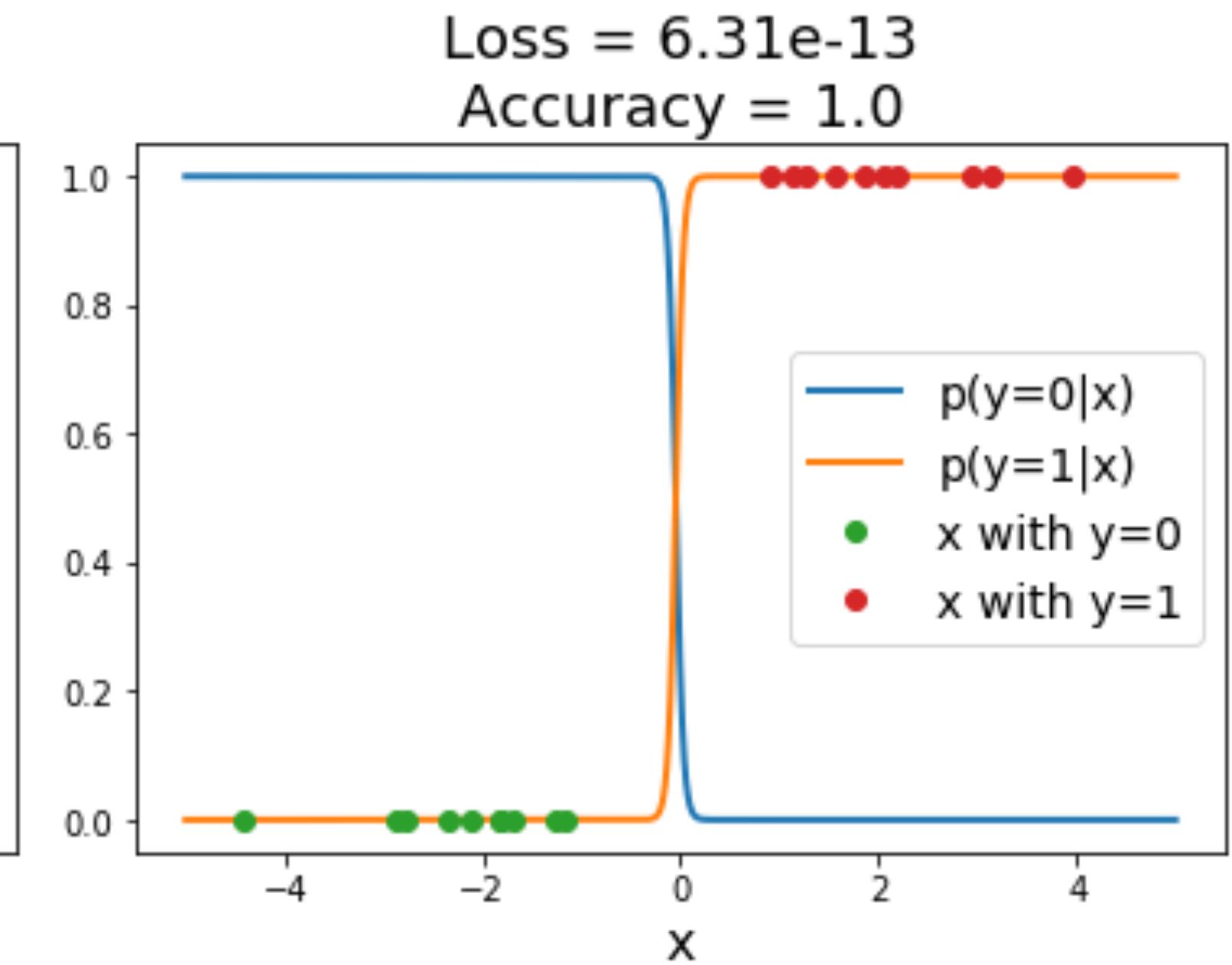
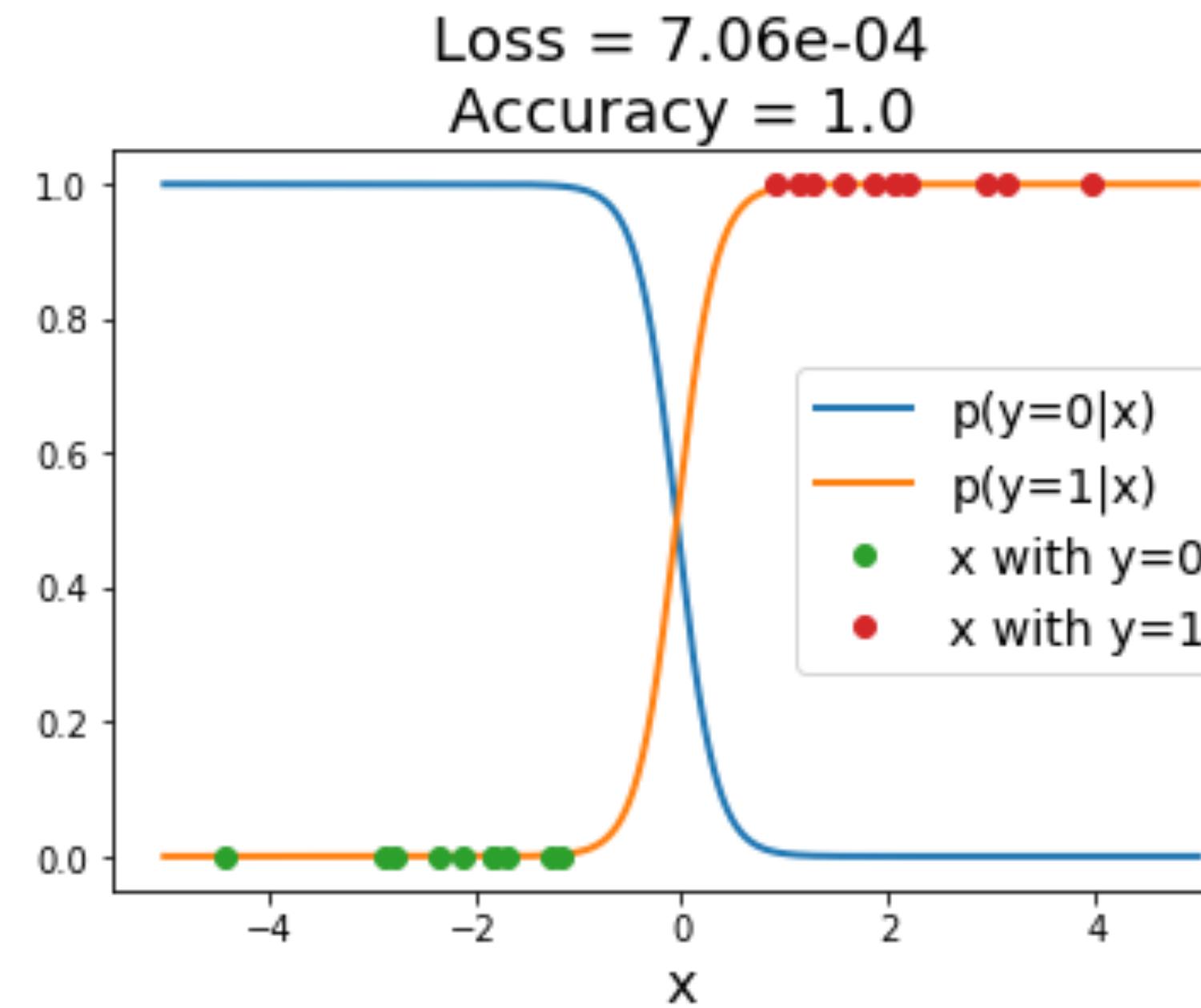
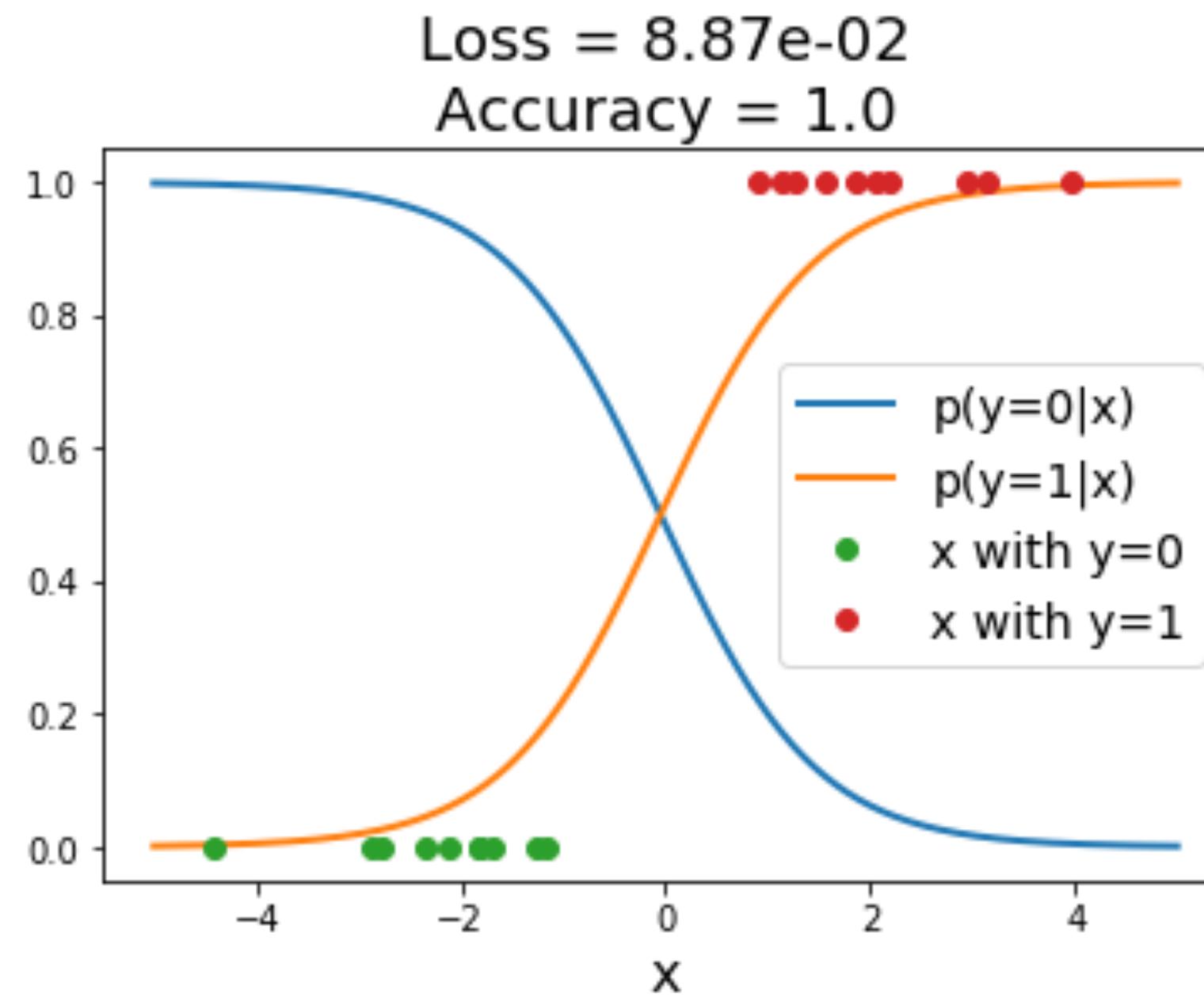
$$L = -\log(p_y)$$



**Both models have perfect accuracy on the training data!**

# Overfitting

A model is overfit when it performs too well on the training data, and has poor performance for unseen data

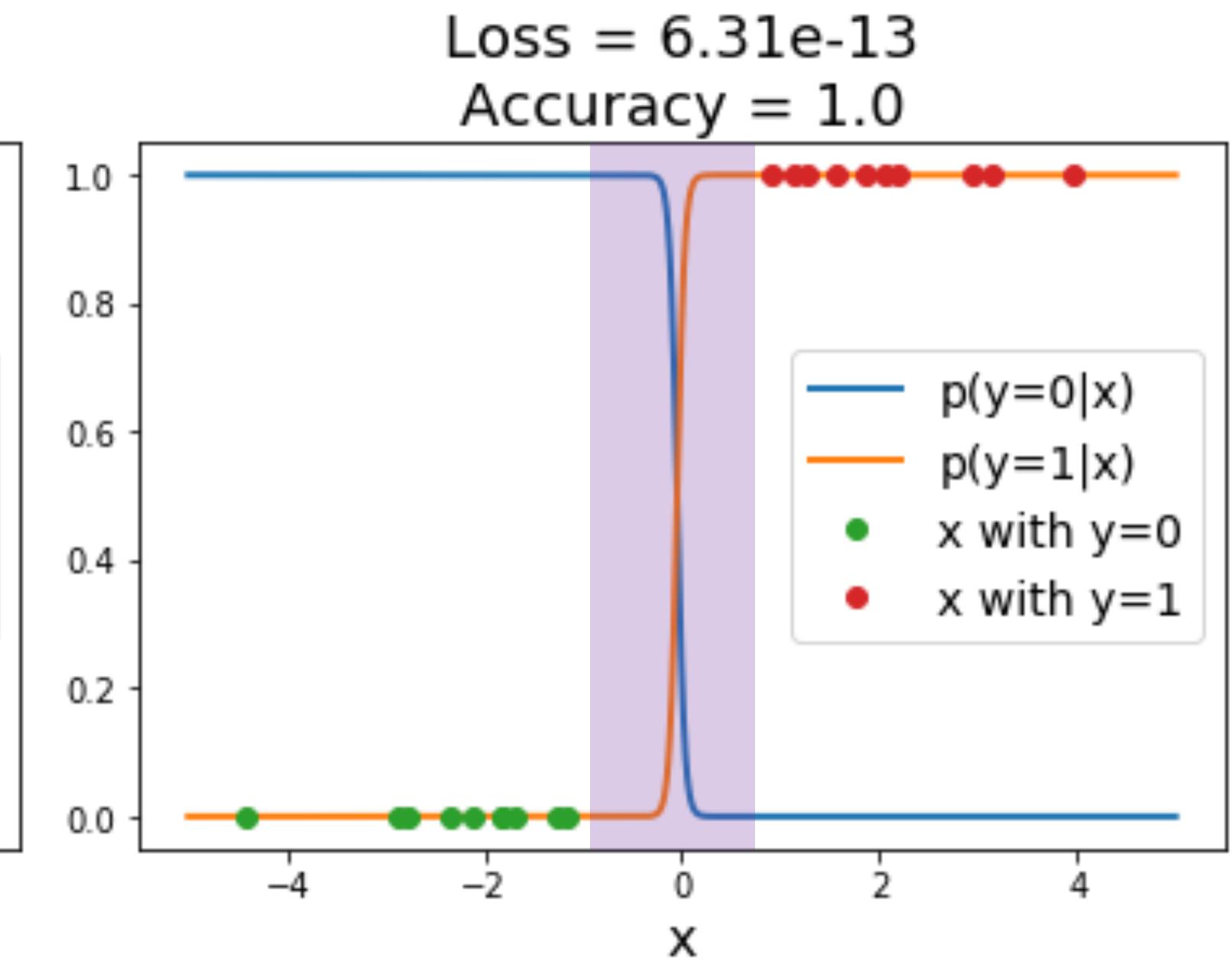
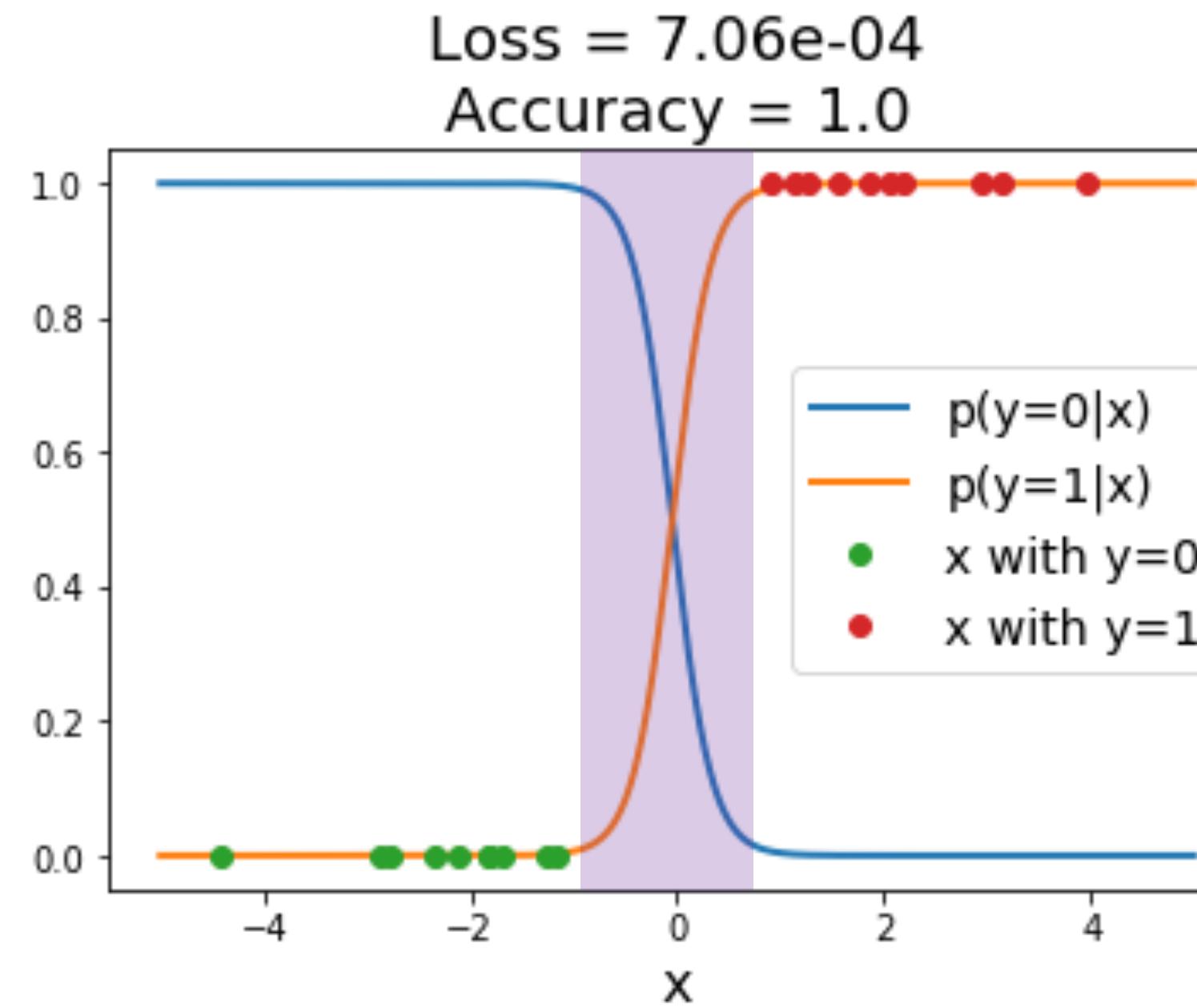
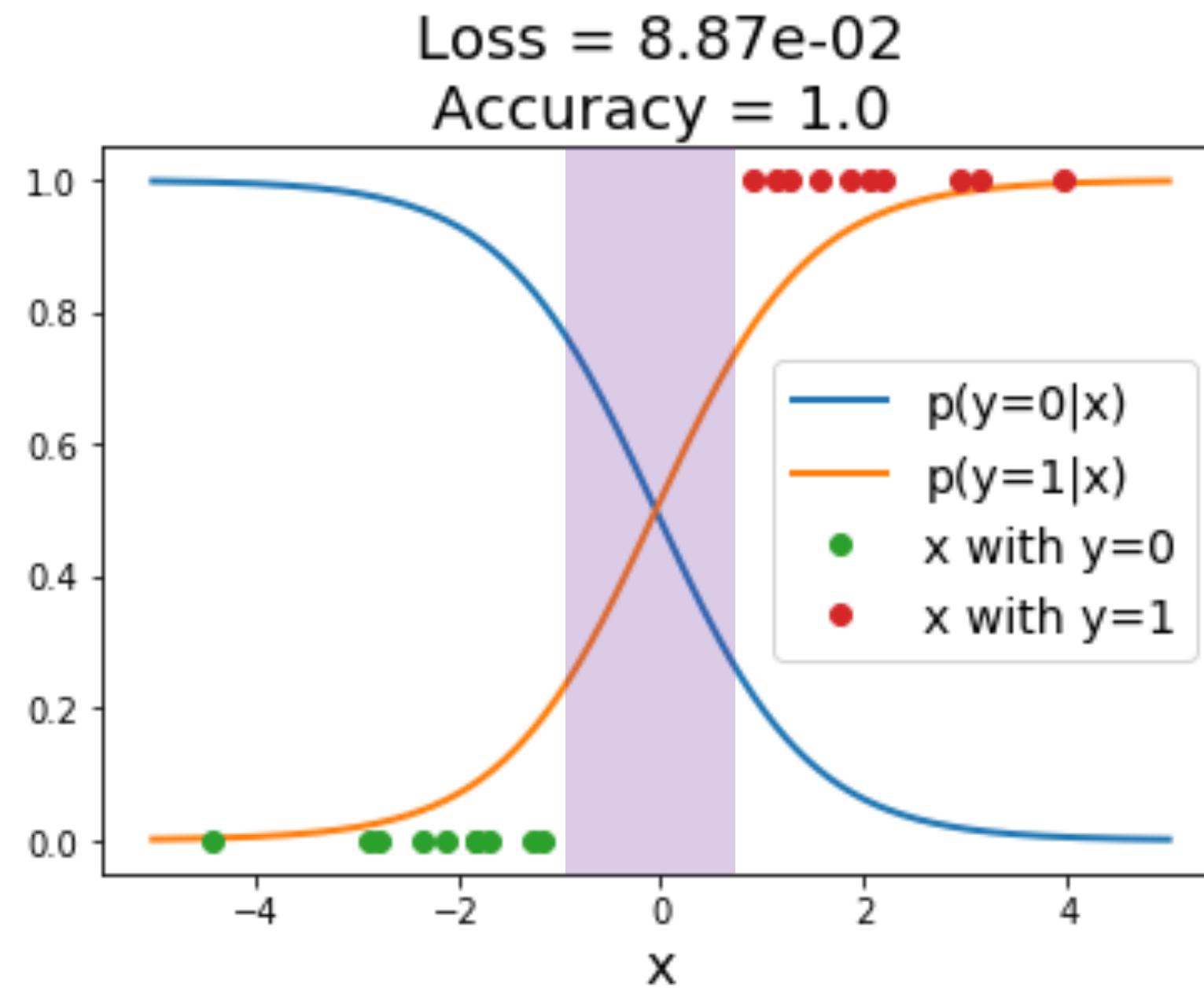


Both models have perfect accuracy on the training data!

Low loss, but unnatural “cliff” between the training points

# Overfitting

A model is overfit when it performs too well on the training data, and has poor performance for unseen data



Overconfidence in regions with no training data could give poor generalization

# Regularization: Beyond Training Error

---

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)$$

**Data loss:** Model predictions  
should match training data



# Regularization: Beyond Training Error

$$L(W) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)}_{\textbf{Data loss}: \text{Model predictions should match training data}} + \underbrace{\lambda R(W)}_{\textbf{Regularization}: \text{Prevent the model from doing too well on training data}}$$

Hyperparameter giving regularization strength

# Regularization: Beyond Training Error

$$L(W) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)}_{\textbf{Data loss}: \text{Model predictions should match training data}} + \underbrace{\lambda R(W)}_{\textbf{Regularization}}$$

Hyperparameter giving regularization strength

**Regularization:** Prevent the model from doing too well on training data

## Simple examples:

L2 regularization:  $R(W) = \sum_{k,l} W_{k,l}^2$

L1 regularization:  $R(W) = \sum_{k,l} |W_{k,l}|$



# Regularization: Beyond Training Error

$$L(W) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)}_{\text{Data loss: Model predictions should match training data}} + \underbrace{\lambda R(W)}_{\text{Regularization}}$$

Hyperparameter giving regularization strength

**Regularization:** Prevent the model from doing too well on training data

## Simple examples:

L2 regularization:  $R(W) = \sum_{k,l} W_{k,l}^2$

L1 regularization:  $R(W) = \sum_{k,l} |W_{k,l}|$

## More complex:

Dropout

Batch normalization

Cutout, Mixup, Stochastic depth, etc...



# Regularization: Prefer Simpler Models

---

Example: Linear classifier with 1D inputs, 2 classes, and softmax loss

$$s_i = w_i x + b_i$$
$$p_i = \frac{\exp(s_i)}{\exp(s_1) + \exp(s_2)}$$
$$L = -\log(p_y) + \lambda \sum_i w_i^2$$

# Regularization: Prefer Simpler Models

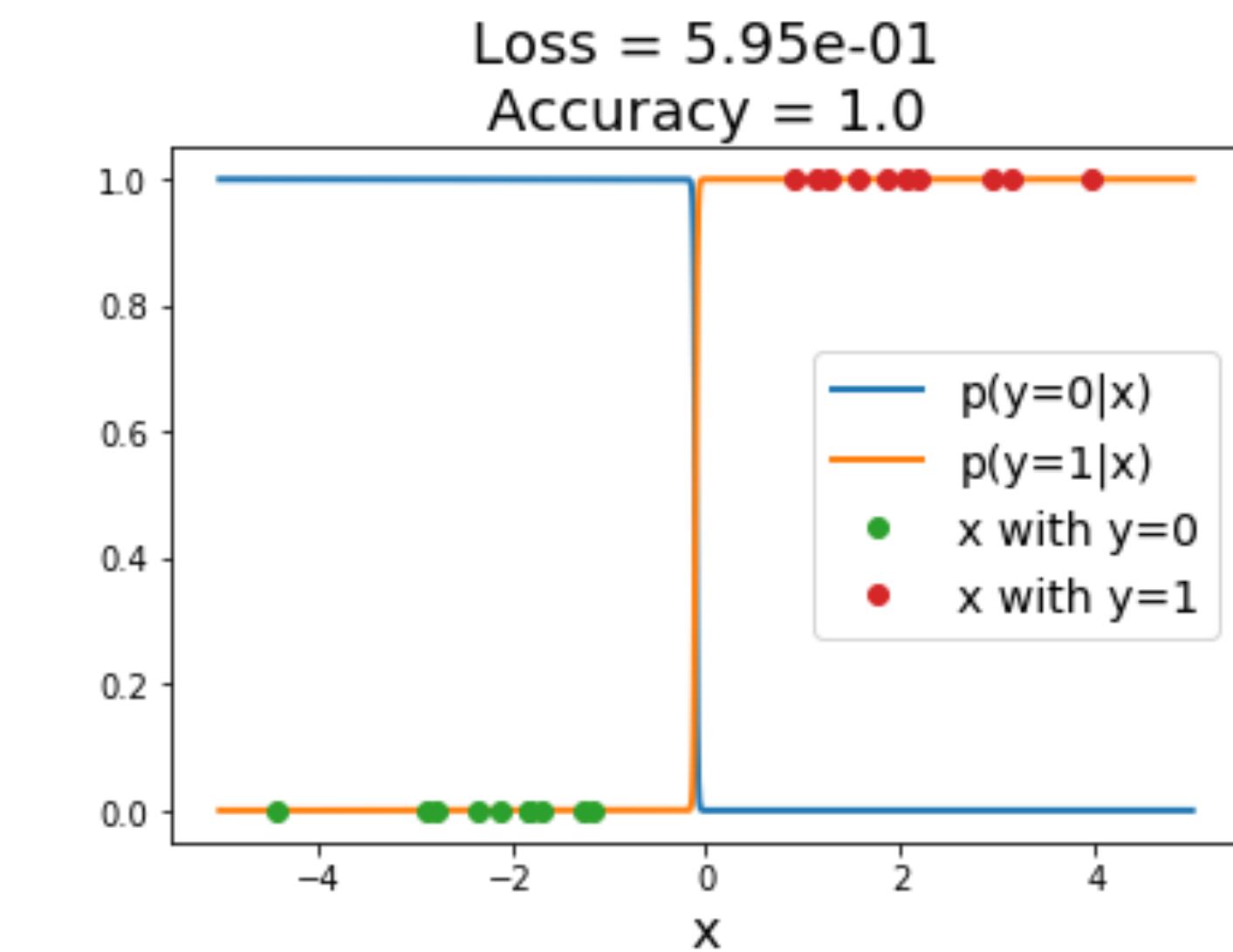
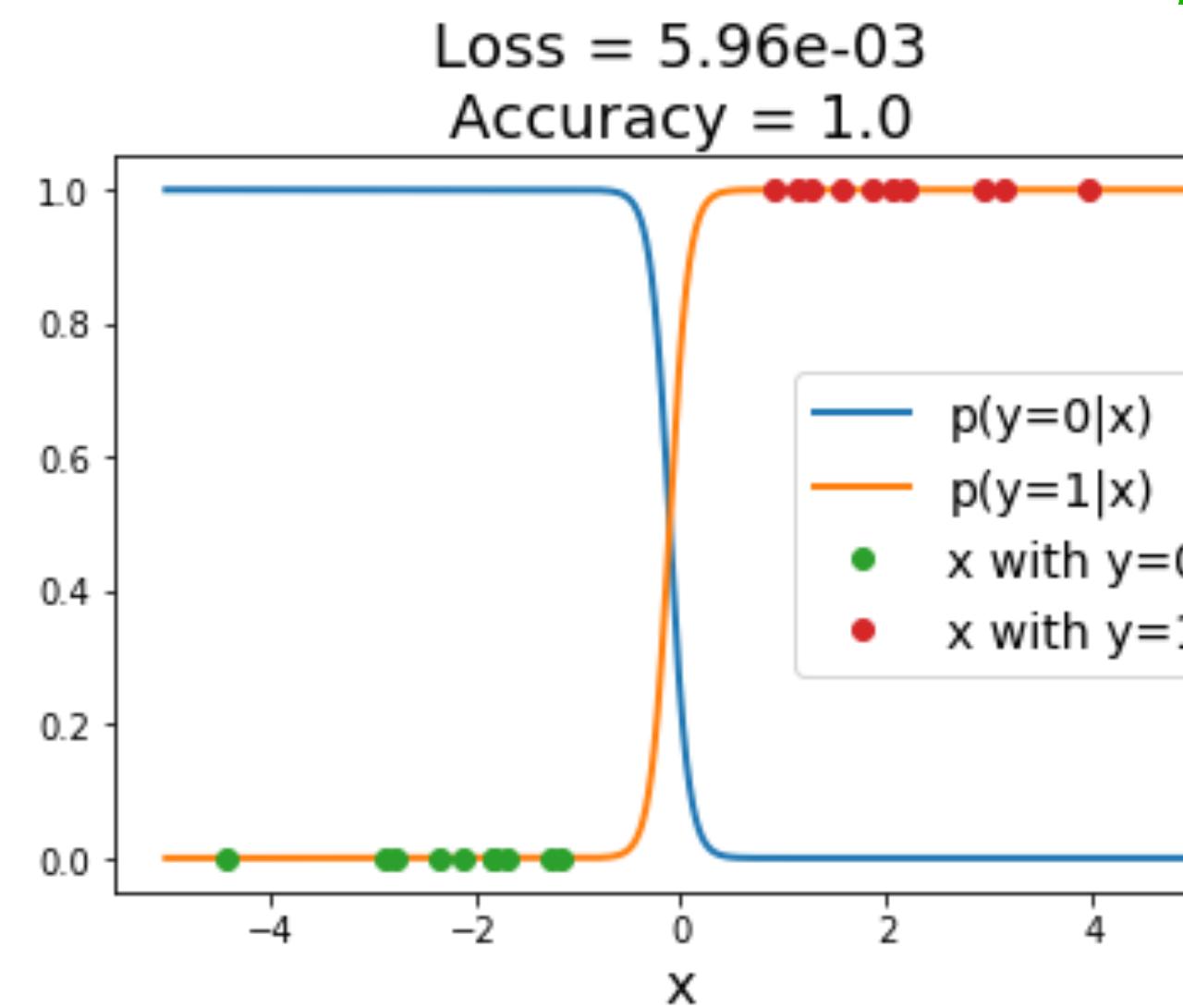
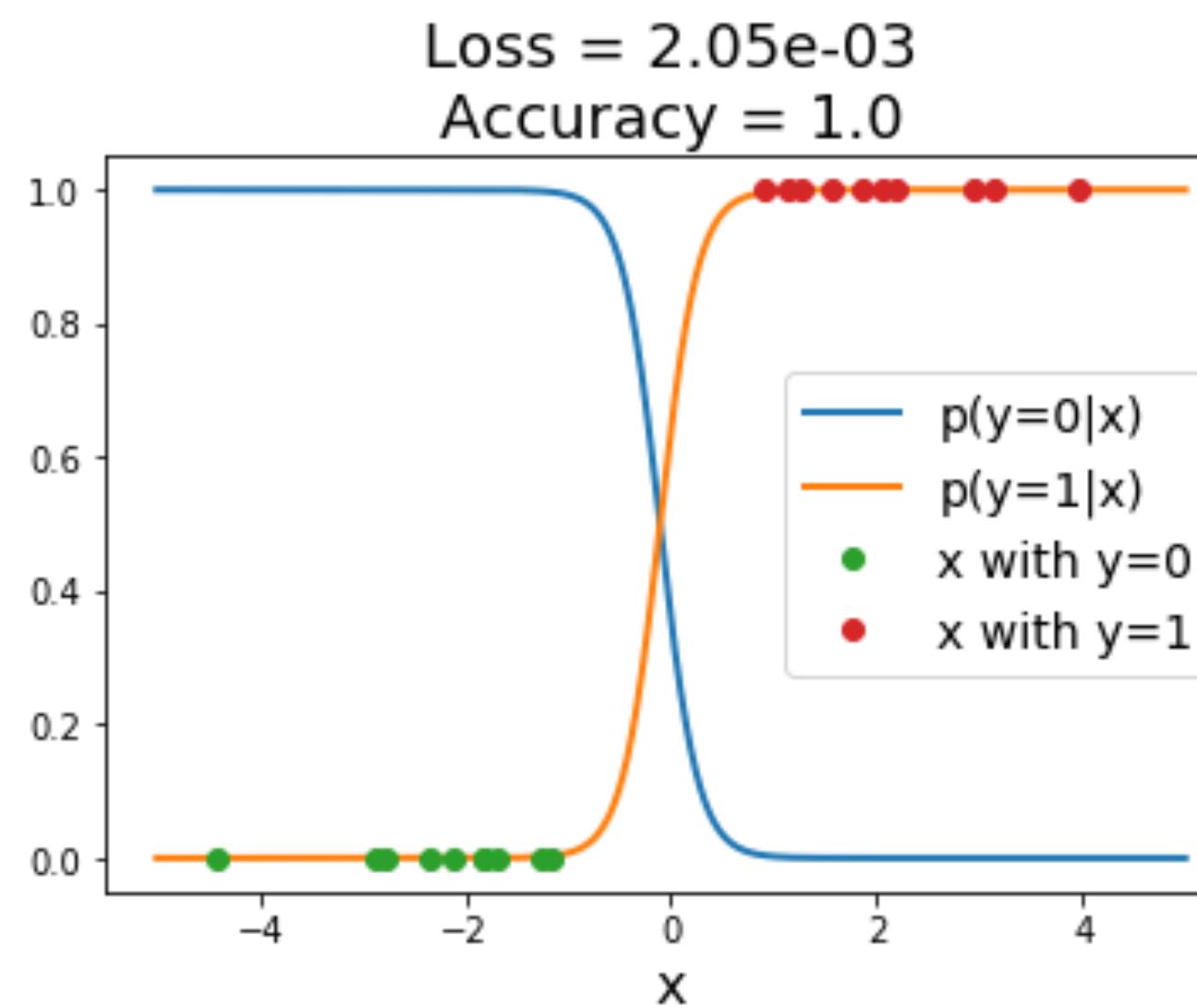
Example: Linear classifier with 1D inputs, 2 classes, and softmax loss

$$s_i = w_i x + b_i$$

$$p_i = \frac{\exp(s_i)}{\exp(s_1) + \exp(s_2)}$$

$$L = -\log(p_y) + \lambda \sum_i w_i^2$$

Regularization term causes loss to **increase** for model with sharp cliff



# Regularization: Expressing Preferences

---

$$x = [1, 1, 1, 1]$$

$$w_1 = [1, 0, 0, 0]$$

$$w_2 = [0.25, 0.25, 0.25, 0.25]$$

L2 Regularization

$$R(W) = \sum_{k,l} W_{k,l}^2$$

$$w_1^T x = w_2^T x = 1$$

Same predictions, so data loss will always be the same

# Regularization: Expressing Preferences

$$x = [1, 1, 1, 1]$$

$$w_1 = [1, 0, 0, 0]$$

$$w_2 = [0.25, 0.25, 0.25, 0.25]$$

L2 Regularization

$$R(W) = \sum_{k,l} W_{k,l}^2$$

L2 Regularization prefers weights to be  
“spread out”

$$w_1^T x = w_2^T x = 1$$

Same predictions, so data loss  
will always be the same



# Finding a good W

---

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i) + \lambda R(W)$$

**Loss function** consists of **data loss** to fit the training data and **regularization** to prevent overfitting

# Optimization

---

$$w^* = \arg \min_w L(w)$$

DR



The [valley image](#) and the [walking man image](#) are in CC0 1.0 public domain



# Idea #1: Random Search (bad idea!)

```
# assume X_train is the data where each column is an example (e.g. 3073 x 50,000)
# assume Y_train are the labels (e.g. 1D array of 50,000)
# assume the function L evaluates the loss function

bestloss = float("inf") # Python assigns the highest possible float value
for num in xrange(1000):
    W = np.random.randn(10, 3073) * 0.0001 # generate random parameters
    loss = L(X_train, Y_train, W) # get the loss over the entire training set
    if loss < bestloss: # keep track of the best solution
        bestloss = loss
        bestW = W
    print 'in attempt %d the loss was %f, best %f' % (num, loss, bestloss)

# prints:
# in attempt 0 the loss was 9.401632, best 9.401632
# in attempt 1 the loss was 8.959668, best 8.959668
# in attempt 2 the loss was 9.044034, best 8.959668
# in attempt 3 the loss was 9.278948, best 8.959668
# in attempt 4 the loss was 8.857370, best 8.857370
# in attempt 5 the loss was 8.943151, best 8.857370
# in attempt 6 the loss was 8.605604, best 8.605604
# ... (truncated: continues for 1000 lines)
```



# Idea #1: Random Search (bad idea!)

```
# Assume X_test is [3073 x 10000], Y_test [10000 x 1]
scores = Wbest.dot(Xte_cols) # 10 x 10000, the class scores for all test examples
# find the index with max score in each column (the predicted class)
Yte_predict = np.argmax(scores, axis = 0)
# and calculate accuracy (fraction of predictions that are correct)
np.mean(Yte_predict == Yte)
# returns 0.1555
```

15.5 % accuracy on CIFAR-10! not bad!



# Idea #1: Random Search (bad idea!)

```
# Assume X_test is [3073 x 10000], Y_test [10000 x 1]
scores = Wbest.dot(Xte_cols) # 10 x 10000, the class scores for all test examples
# find the index with max score in each column (the predicted class)
Yte_predict = np.argmax(scores, axis = 0)
# and calculate accuracy (fraction of predictions that are correct)
np.mean(Yte_predict == Yte)
# returns 0.1555
```

15.5 % accuracy on CIFAR-10! not bad!  
(SOTA is ~95%)



# Idea #2: Follow the slope

---



# Idea #2: Follow the slope

---

In 1-dimension, the **derivative** of a function gives the slope:

$$\frac{df}{dx} = \lim_{h \rightarrow 0} \frac{f(x + h) - f(x)}{h}$$



# Idea #2: Follow the slope

---

In 1-dimension, the **derivative** of a function gives the slope:

$$\frac{df}{dx} = \lim_{h \rightarrow 0} \frac{f(x + h) - f(x)}{h}$$

In multiple dimensions, the **gradient** is the vector of (partial derivatives) along each dimension

The slope in any direction is the **dot product** of the direction with the gradient. The direction of steepest descent is the **negative gradient**.





Current W:

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33, ...]

loss 1.25347

Gradient  $\frac{dL}{dW}$   
[?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?, ...,]





Current W:

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33, ...]

loss 1.25347

W + h (first dim):

[0.34 + 0.0001,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33, ...]

loss 1.25322

Gradient  $\frac{dL}{dW}$

[?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?, ...,]



Current  $\mathbf{W}$ :

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33, ...]

loss 1.25347

$\mathbf{W} + \mathbf{h}$  (first dim):

[0.34 + 0.0001,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33, ...]

loss 1.25322

Gradient  $\frac{dL}{dW}$

[-2.5,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
...,  
?]

$$\frac{(1.25322 - 1.25347)}{0.0001} = -2.5$$

$$\frac{df}{dx} = \lim_{h \rightarrow 0} \frac{f(x + h) - f(x)}{h}$$

?, ...]

Current W:

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33, ...]

loss 1.25347

W + h (second dim):

[0.34,  
-1.11 + **0.0001**,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33, ...]

loss 1.25353

Gradient  $\frac{dL}{dW}$

[-2.5,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?, ..., ]

Current W:

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33, ...]

loss 1.25347

W + h (second dim):

[0.34,  
-1.11 + **0.0001**,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33, ...]

loss 1.25353

Gradient  $\frac{dL}{dW}$

[-2.5,  
**0.6**,  
?,  
?]

$$\frac{(1.25353 - 1.25347)}{0.0001} = 0.6$$

$$\frac{df}{dx} = \lim_{h \rightarrow 0} \frac{f(x + h) - f(x)}{h}$$



Current  $\mathbf{W}$ :

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33, ...]  
**loss 1.25347**

$\mathbf{W} + \mathbf{h}$  (third dim):

[0.34,  
-1.11,  
0.78 + **0.0001**,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33, ...]  
**loss 1.25353**

Gradient  $\frac{dL}{dW}$

[-2.5,  
0.6,  
**0.0**,  
?,  
?]

$$\frac{(1.25347 - 1.25347)}{0.0001} = 0.0$$

$$\frac{df}{dx} = \lim_{h \rightarrow 0} \frac{f(x + h) - f(x)}{h}$$



Current W:

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33, ...]

loss 1.25347

W + h (third dim):

[0.34,  
-1.11,  
0.78 + 0.0001,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33, ...]

loss 1.25353

Gradient  $\frac{dL}{dW}$

[-2.5,  
0.6,  
0.0,  
?,  
?,  
?,  
?,  
?,  
?,  
?]

**Numeric Gradient:**

- Slow: O(#dimensions)
- Approximate

# Loss is a function of W

---

$$L = \frac{1}{2} \sum_{i=1}^N L_i + \sum_k W_k^2$$

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$s = f(x, W) = Wx$$

Want  $\nabla_w L$



# Loss is a function of W

---

$$L = \frac{1}{2} \sum_{i=1}^N L_i + \sum_k W_k^2$$

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$s = f(x, W) = Wx$$

Use calculus to compute an  
**analytic gradient**

Want  $\nabla_w L$

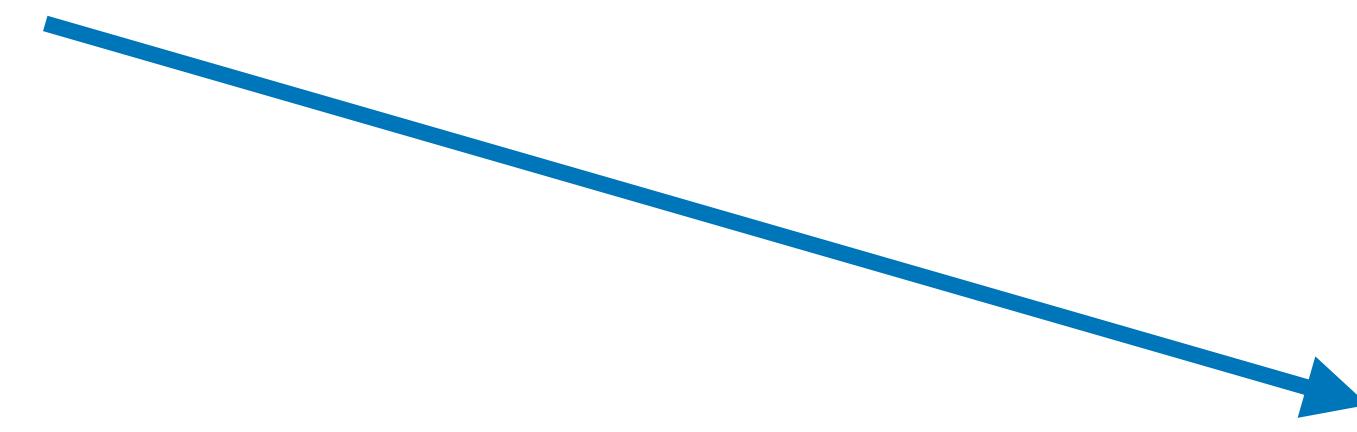


Current W:

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33, ...]

loss 1.25347

$\frac{dL}{dW}$  = some function of data and  $W$



Gradient  $\frac{dL}{dW}$   
[-2.5,  
0.6,  
0.0,  
0.2,  
0.7,  
-0.5,  
1.1,  
1.3,  
-2.1, ...]



Current W:

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33, ...]

loss 1.25347

Gradient  $\frac{dL}{dW}$

[-2.5,  
0.6,  
0.0,  
0.2,  
0.7,  
-0.5,  
1.1,  
1.3,  
-2.1, ...]

$\frac{dL}{dW}$  = some function of data and  $W$

In practice we will compute  $\frac{dL}{dW}$   
using back propagation;  
see Lecture 6

# Computing Gradients

---

- **Numeric gradient:** approximate, slow, easy to write
- **Analytic gradient:** exact, fast, error-prone

# Computing Gradients

---

- **Numeric gradient:** approximate, slow, easy to write
- **Analytic gradient:** exact, fast, error-prone

In practice: Always use analytic gradient, but check implementation with numerical gradient. This is called a **gradient check**.

# Computing Gradients

- **Numeric gradient:** approximate, slow, easy to write
- **Analytic gradient:** exact, fast, error-prone

In practice: Always use analytic gradient, but check implementation with numerical gradient. This is called a **gradient check**.

```
def grad_check_sparse(f, x, analytic_grad, num_checks=10, h=1e-7):
    """
    sample a few random elements and only return numerical
    in this dimensions.
    """

```

# Computing Gradients

- **Numeric gradient:** approximate, slow, easy to write
- **Analytic gradient:** exact, fast, error-prone

```
torch.autograd.gradcheck(func, inputs, eps=1e-06, atol=1e-05, rtol=0.001,  
raise_exception=True, check_sparse_nnz=False, nondet_tol=0.0)
```

[SOURCE] 

Check gradients computed via small finite differences against analytical gradients w.r.t. tensors in `inputs` that are of floating point type and with `requires_grad=True`.

The check between numerical and analytical gradients uses `allclose()`.



# Computing Gradients

- **Numeric gradient:** approximate, slow, easy to write
- **Analytic gradient:** exact, fast, error-prone

```
torch.autograd.gradgradcheck(func, inputs, grad_outputs=None, eps=1e-06, atol=1e-05, rtol=0.001, gen_non_contig_grad_outputs=False, raise_exception=True, nondet_tol=0.0)
```

[\[SOURCE\]](#)

Check gradients of gradients computed via small finite differences against analytical gradients w.r.t. tensors in `inputs` and `grad_outputs` that are of floating point type and with `requires_grad=True`.

This function checks that backpropagating through the gradients computed to the given `grad_outputs` are correct.



# Gradient Descent

---

- Iteratively step in the direction of the negative gradient (direction of local steepest descent)

```
# Vanilla gradient descent
w = initialize_weights()
for t in range(num_steps):
    dw = compute_gradient(loss_fn, data, w)
    w -= learning_rate * dw
```

## Hyperparameters:

- Weight initialization method
- Number of steps
- Learning rate



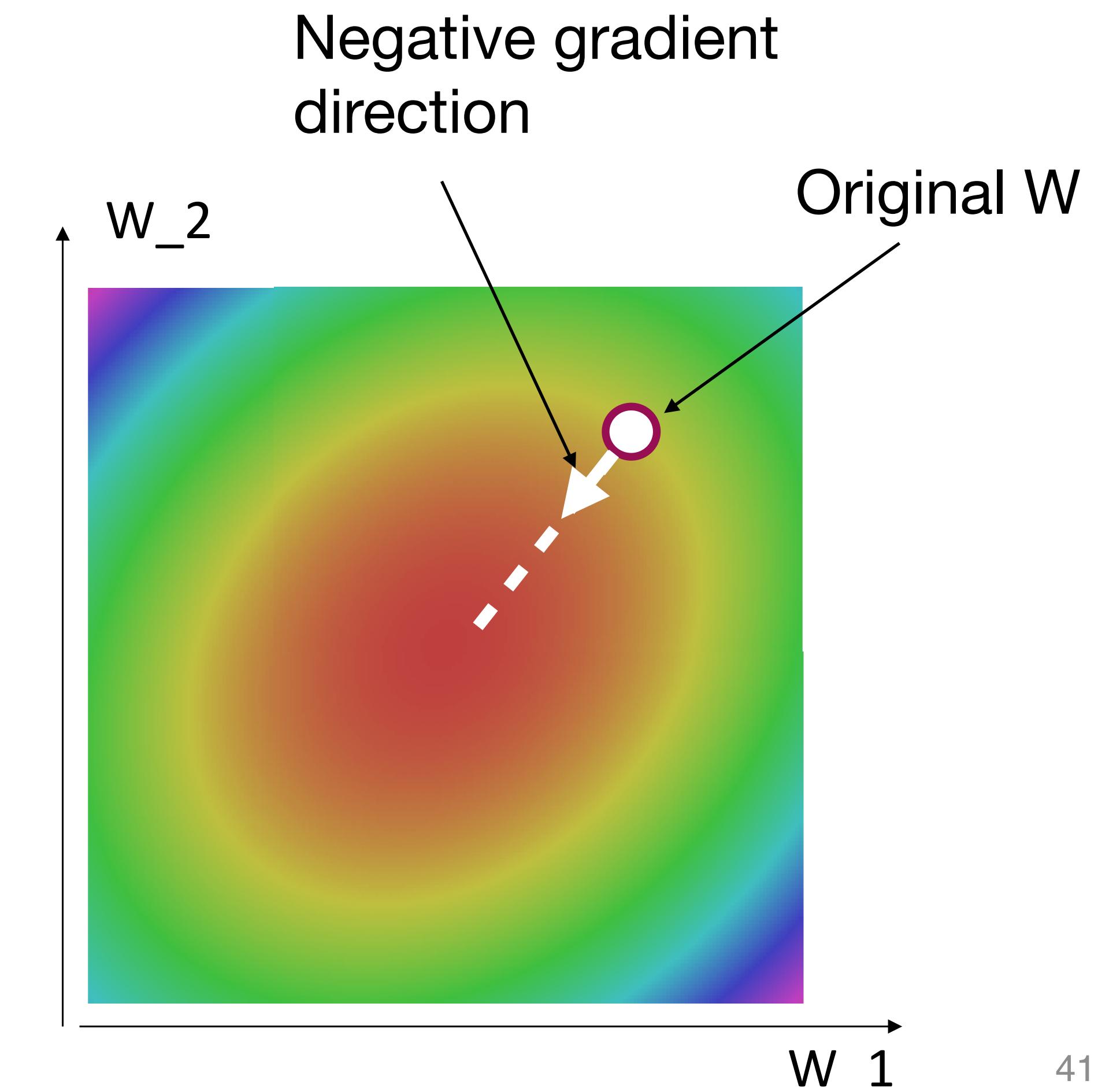
# Gradient Descent

- Iteratively step in the direction of the negative gradient (direction of local steepest descent)

```
# Vanilla gradient descent
w = initialize_weights()
for t in range(num_steps):
    dw = compute_gradient(loss_fn, data, w)
    w -= learning_rate * dw
```

## Hyperparameters:

- Weight initialization method
- Number of steps
- Learning rate



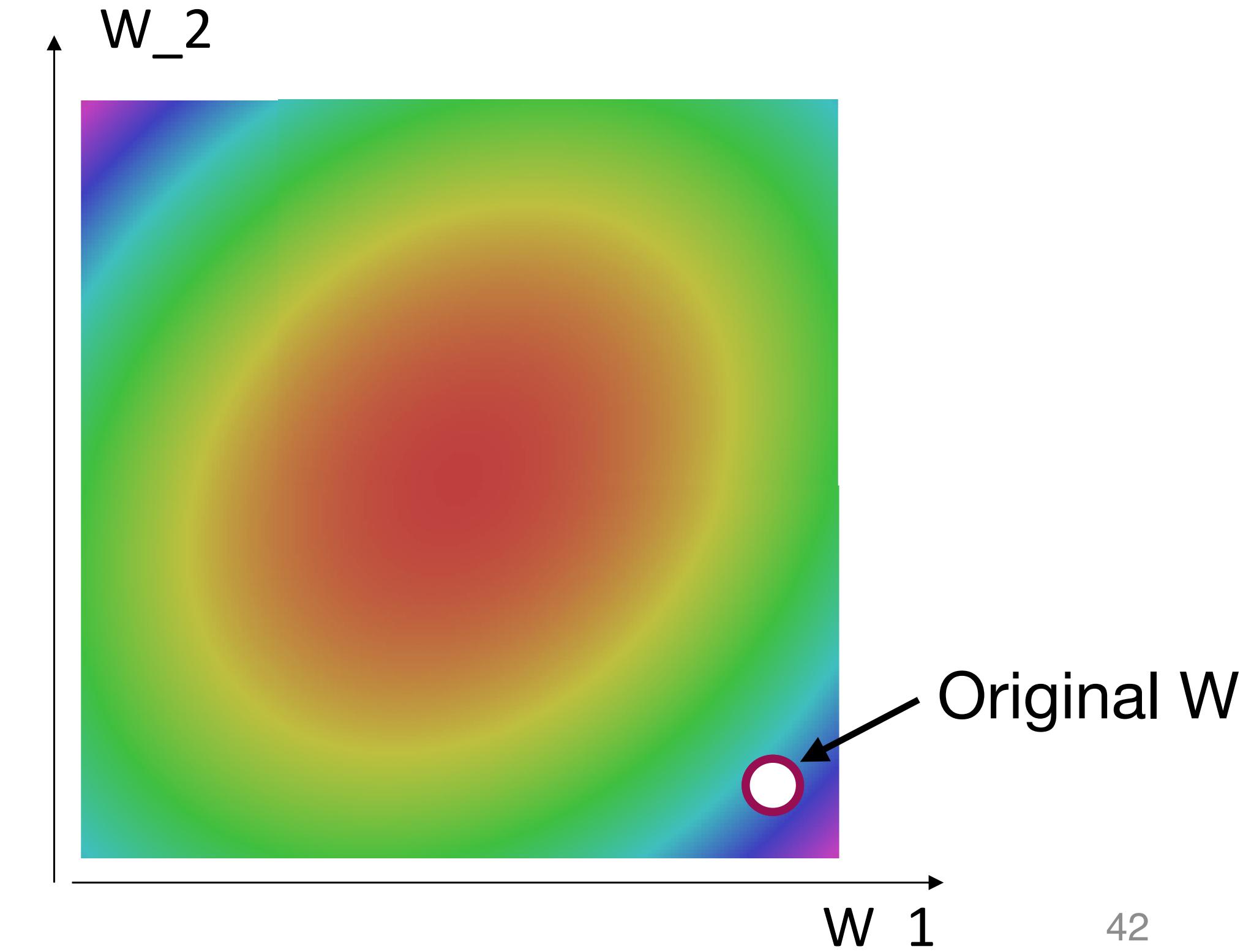
# Gradient Descent

- Iteratively step in the direction of the negative gradient (direction of local steepest descent)

```
# Vanilla gradient descent
w = initialize_weights()
for t in range(num_steps):
    dw = compute_gradient(loss_fn, data, w)
    w -= learning_rate * dw
```

## Hyperparameters:

- Weight initialization method
- Number of steps
- Learning rate



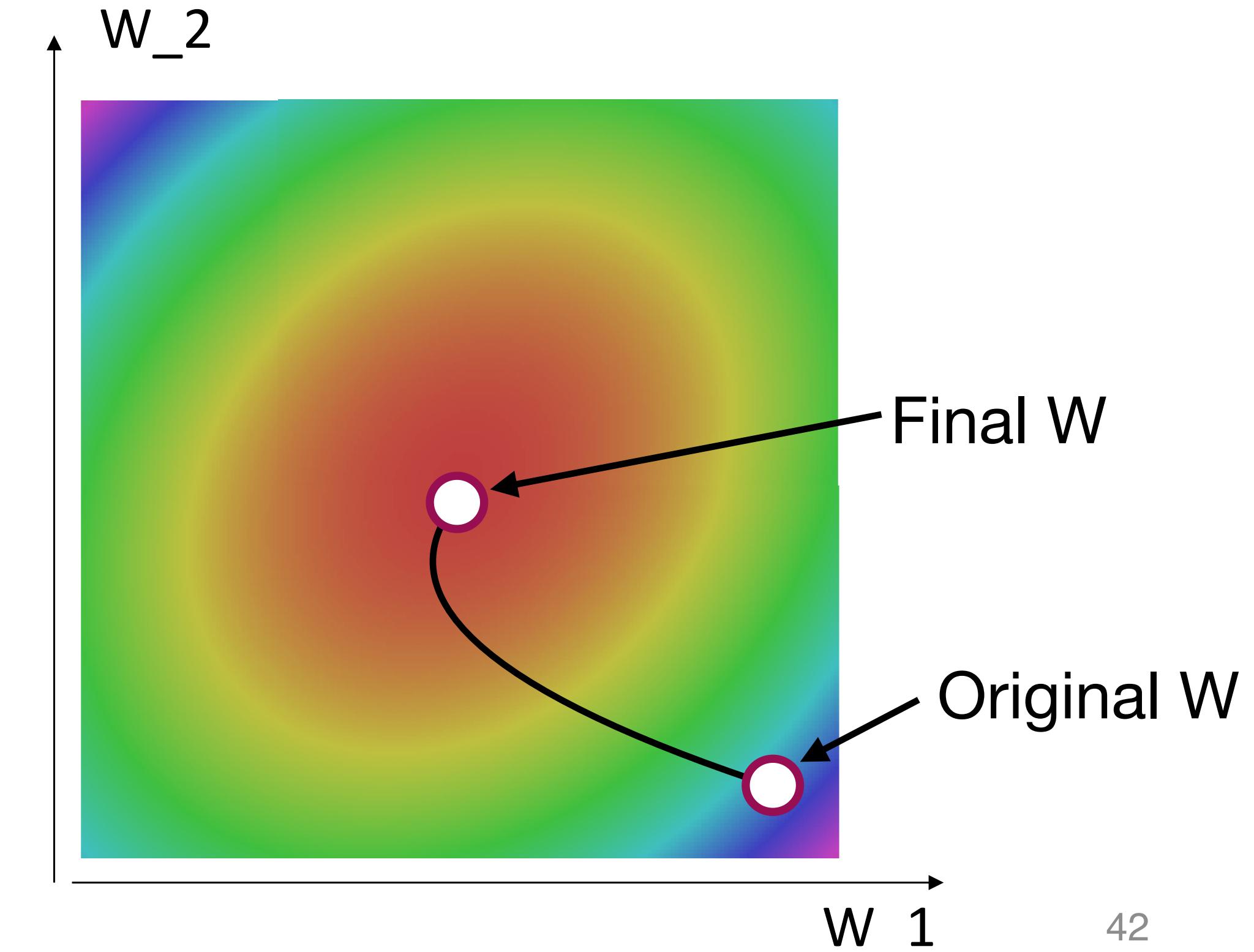
# Gradient Descent

- Iteratively step in the direction of the negative gradient (direction of local steepest descent)

```
# Vanilla gradient descent
w = initialize_weights()
for t in range(num_steps):
    dw = compute_gradient(loss_fn, data, w)
    w -= learning_rate * dw
```

## Hyperparameters:

- Weight initialization method
- Number of steps
- Learning rate



# Batch Gradient Descent

---

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W) + \lambda R(W)$$

Full sum expensive  
when N is large!

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$



# Stochastic Gradient Descent (SGD)

---

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W) + \lambda R(W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

```
# Stochastic gradient descent
w = initialize_weights()
for t in range(num_steps):
    minibatch = sample_data(data, batch_size)
    dw = compute_gradient(loss_fn, minibatch, w)
    w -= learning_rate * dw
```

Full sum expensive  
when N is large!

Approximate sum using  
**minibatch** of examples  
32/64/128 common

## Hyperparameters:

- Weight initialization
- Number of steps
- Learning rate
- Batch size
- Data sampling



# Stochastic Gradient Descent (SGD)

---

$$L(W) = \mathbb{E}_{(x,y) \sim p_{data}}[L(x, y, W)] + \lambda R(W)$$

$$\approx \frac{1}{N} \sum_{i=1}^N L(x_i, y_i, W) + \lambda R(W)$$

Think of loss as an expectation  
over the full **data distribution**

$p_{data}$

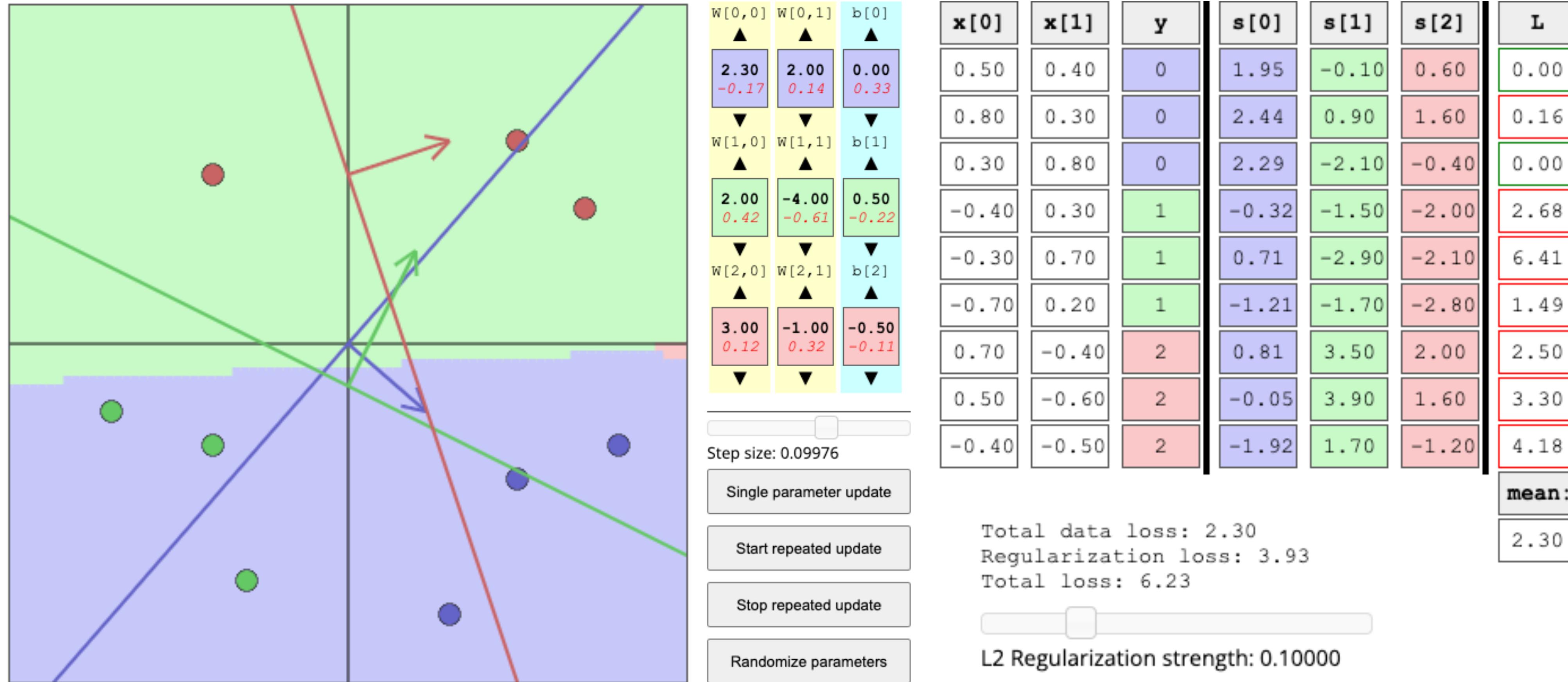
Approximate expectation  
via sampling

$$\nabla_W L(W) = \nabla_W \mathbb{E}_{(x,y) \sim p_{data}}[L(x, y, W)] + \lambda R(W)$$

$$\approx \sum_{i=1}^N N \nabla_w L(x_i, y_i, W) + \nabla_w \lambda R(W)$$



# Interactive Web Demo

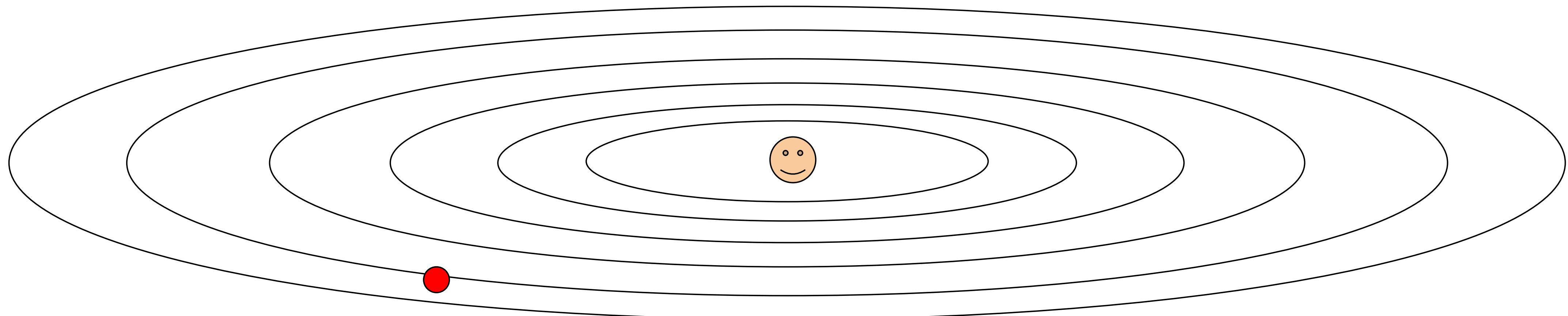


<http://vision.stanford.edu/teaching/cs231n-demos/linear-classify/>

# Problems with SGD

What if loss changes quickly in one direction and slowly in another?

What does gradient decent do?



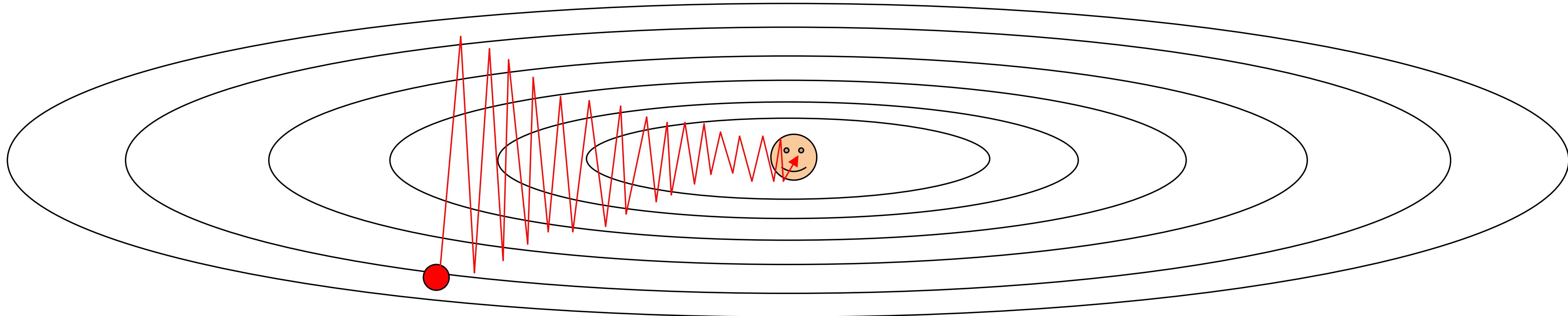
Loss function has high condition number: ratio of largest to smallest singular value of the Hessian matrix is large

# Problems with SGD

What if loss changes quickly in one direction and slowly in another?

What does gradient decent do?

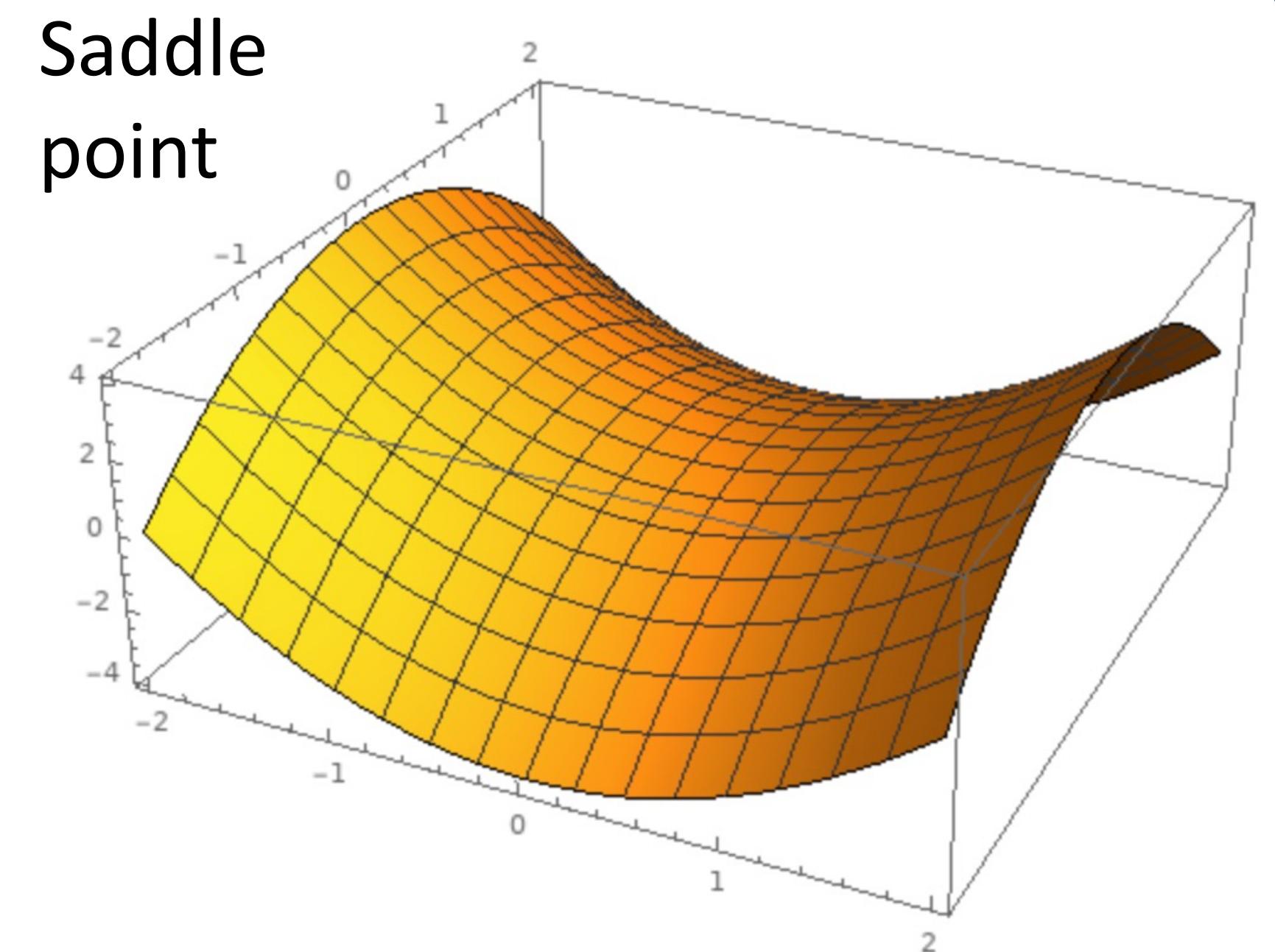
Very slow progress along shallow dimension, jitter along steep direction



Loss function has high condition number: ratio of largest to smallest singular value of the Hessian matrix is large

# Problems with SGD

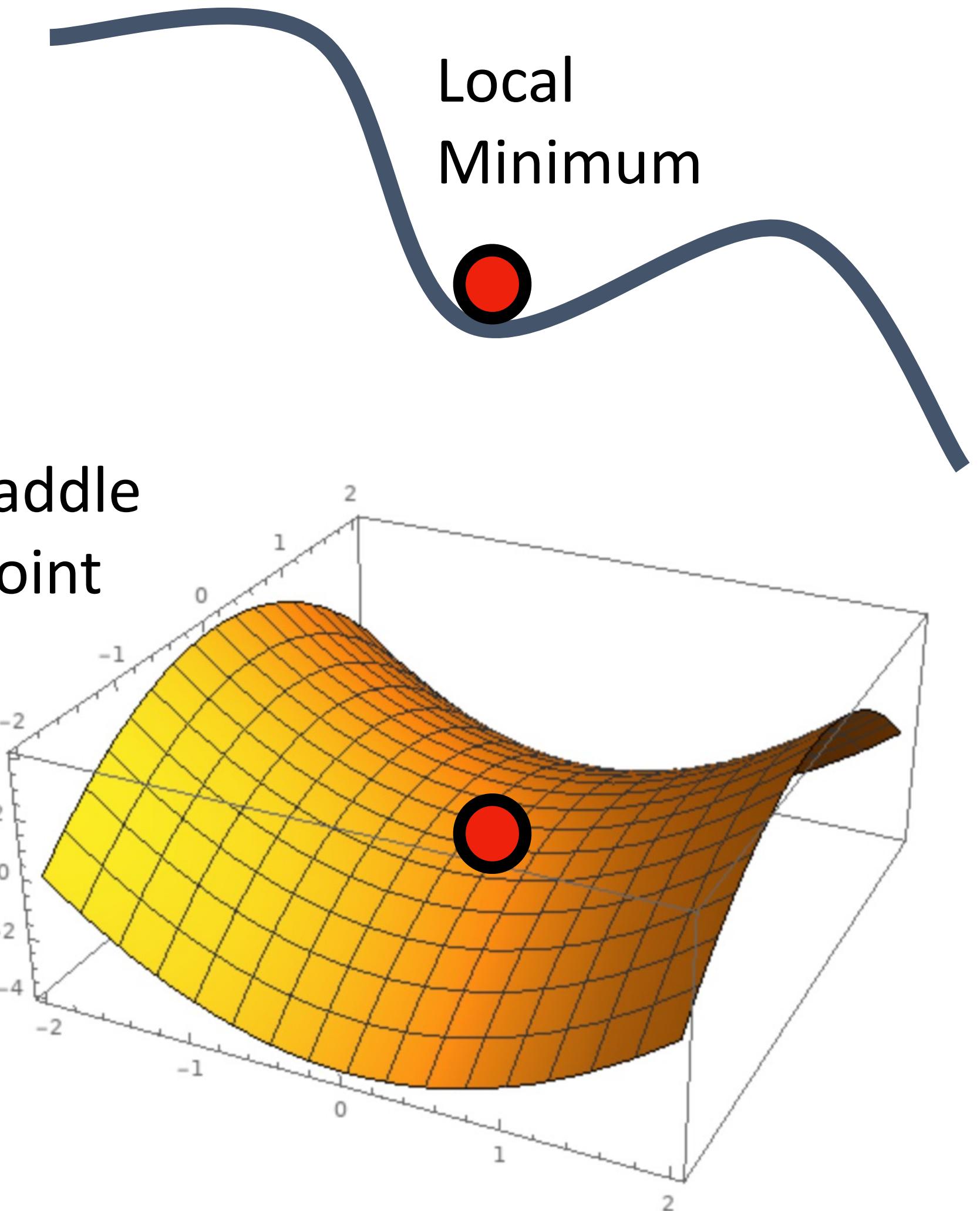
What if the loss function has a  
**local minimum or saddle point?**



# Problems with SGD

What if the loss function has a  
**local minimum or saddle point?**

Zero gradient, gradient descent gets stuck



# Problems with SGD

---

Our gradients come from mini batches so they can be noisy!

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W) + \lambda R(W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

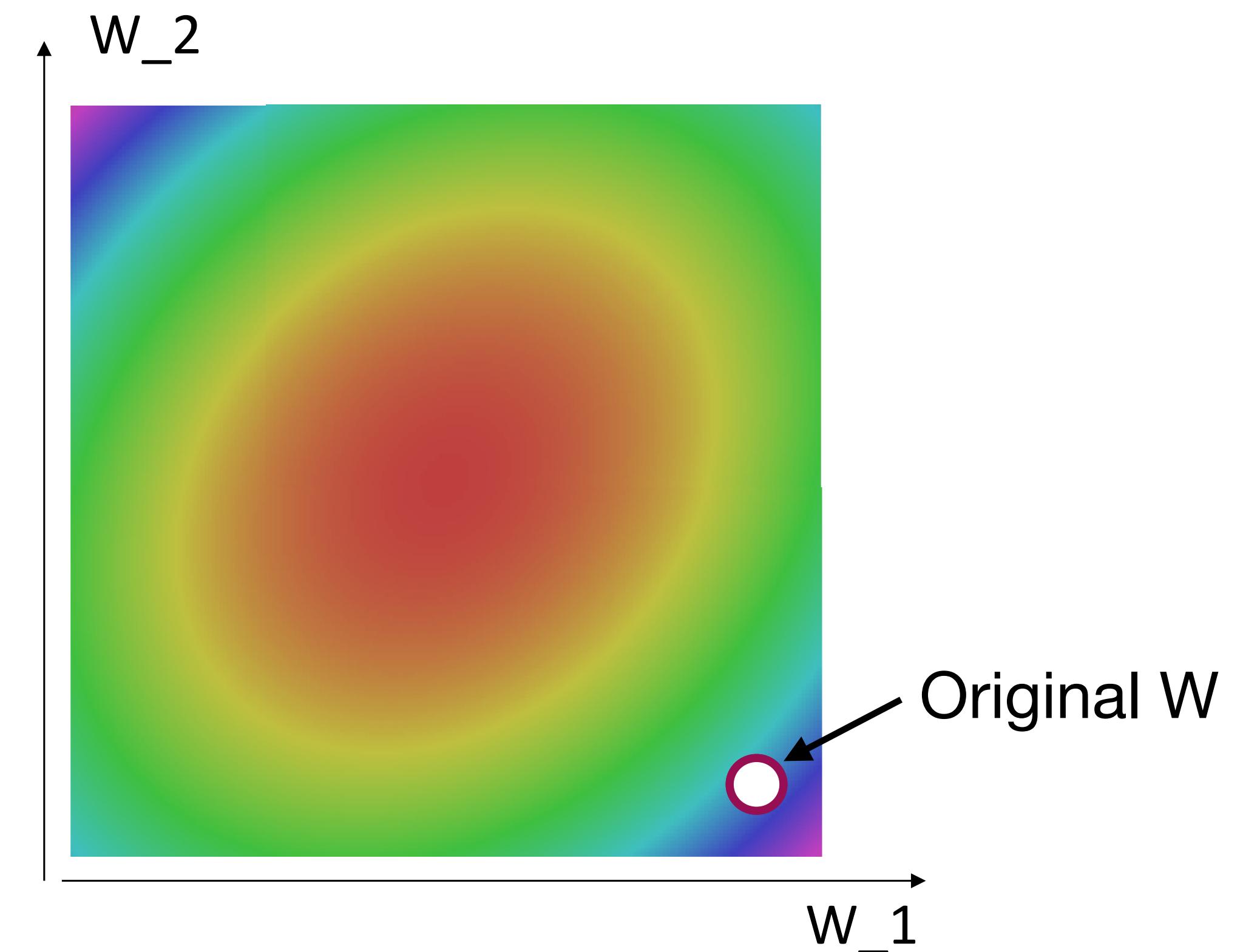


# Problems with SGD

Our gradients come from mini batches so they can be noisy!

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W) + \lambda R(W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

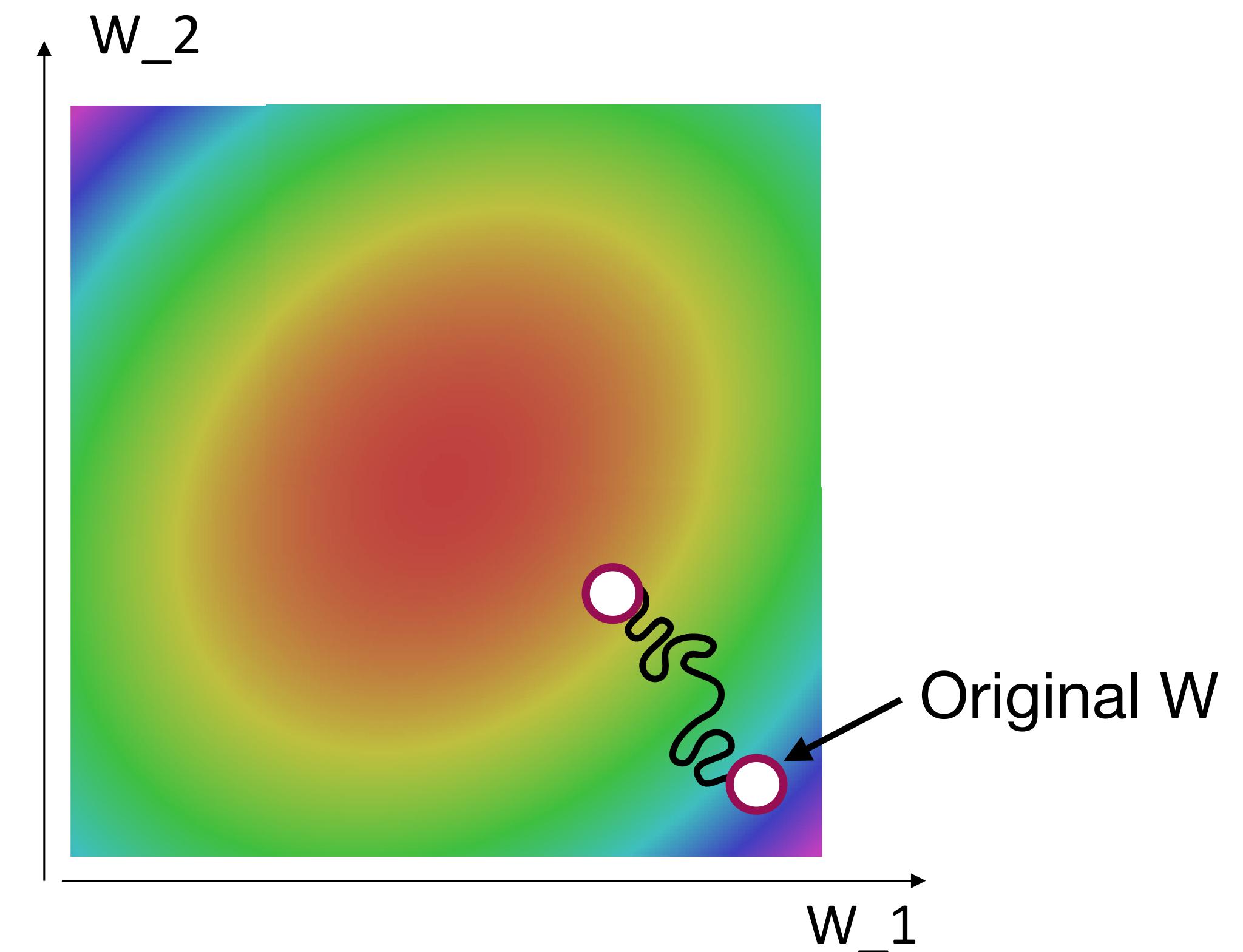


# Problems with SGD

Our gradients come from mini batches so they can be noisy!

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W) + \lambda R(W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

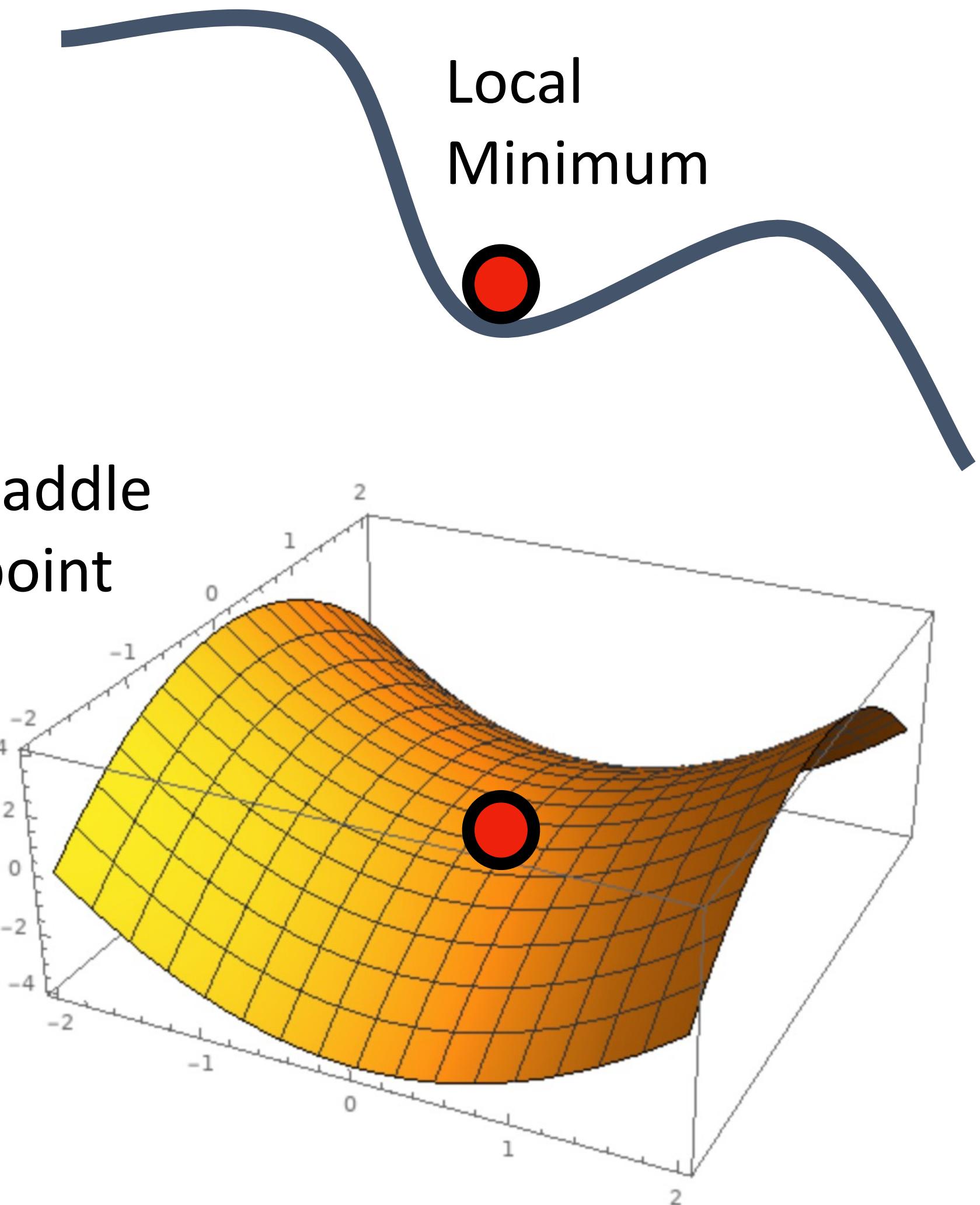


# Problems with SGD

What if the loss function has a **local minimum** or **saddle point**?

Batched gradient descent always computes same gradients

SGD computes noisy gradients, may help to escape saddle points



# SGD + Momentum

---

## SGD

$$w_{t+1} = w_t - \alpha \nabla L(w_t)$$

```
for t in range(num_steps):
    dw = compute_gradient(w)
    w -= learning_rate * dw
```



# SGD + Momentum

SGD

$$w_{t+1} = w_t - \alpha \nabla L(w_t)$$

```
for t in range(num_steps):
    dw = compute_gradient(w)
    w -= learning_rate * dw
```

SGD + Momentum

$$\begin{aligned} v_{t+1} &= \rho v_t + \nabla L(w_t) \\ w_{t+1} &= w_t - \alpha v_{t+1} \end{aligned}$$

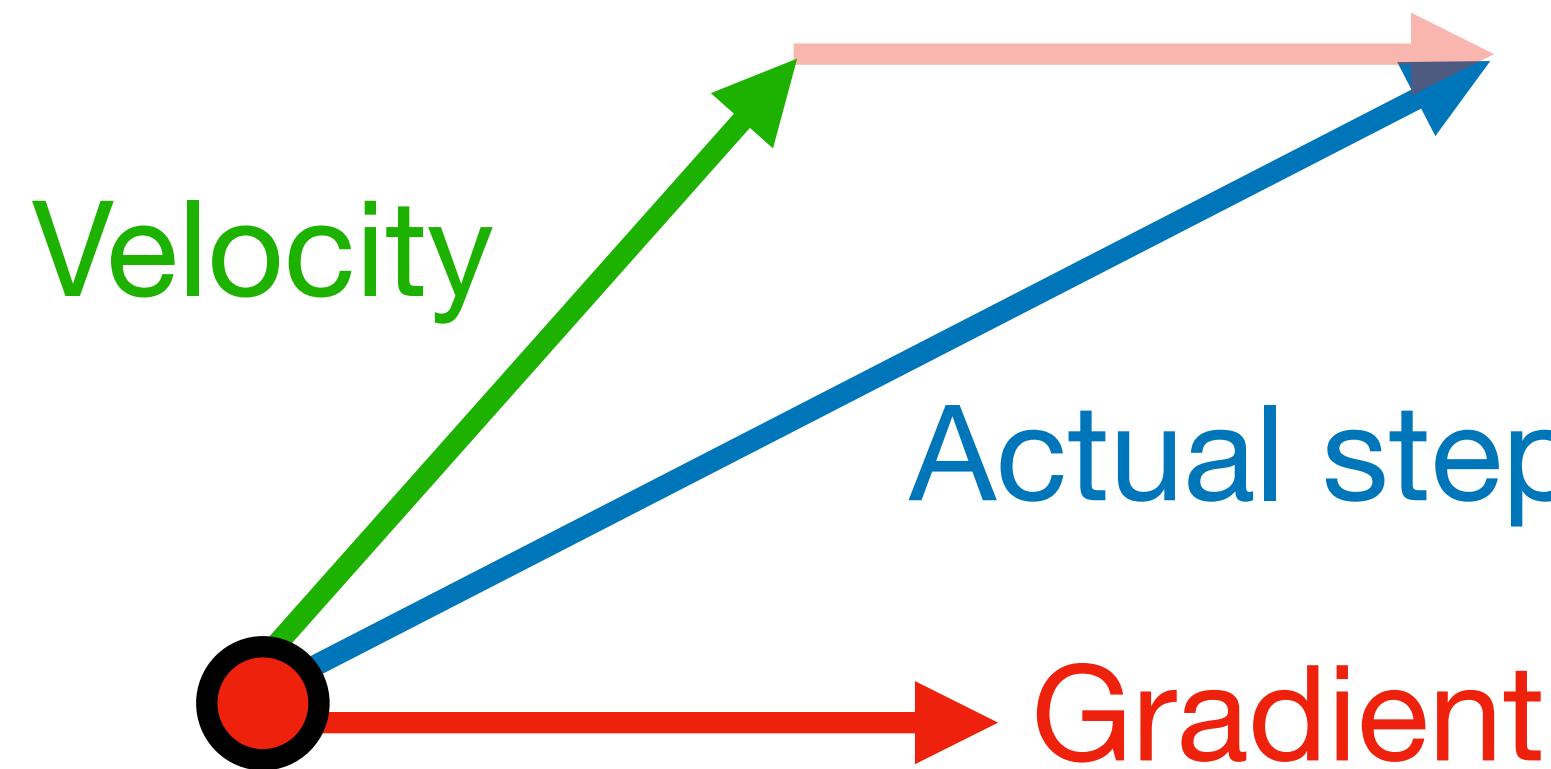
```
v = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    v = rho * v + dw
    w -= learning_rate * v
```

- Build up “velocity” as a running mean of gradients
- Rho gives “friction”; typically rho = 0.9 or 0.99



# SGD + Momentum

Momentum update:



Combine gradient at current point  
with velocity to get step used to  
update weights

SGD + Momentum

$$v_{t+1} = \rho v_t + \nabla L(w_t)$$
$$w_{t+1} = w_t - \alpha v_{t+1}$$

```
v = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    v = rho * v + dw
    w -= learning_rate * v
```

- Build up “velocity” as a running mean of gradients
- Rho gives “friction”; typically rho = 0.9 or 0.99

# SGD + Momentum

---

SGD + Momentum

$$v_{t+1} = \rho v_t - \alpha \nabla L(w_t)$$

$$w_{t+1} = w_t + v_{t+1}$$

```
v = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    v = rho * v - learning_rate * dw
    w += v
```

SGD + Momentum

$$v_{t+1} = \rho v_t + \nabla L(w_t)$$

$$w_{t+1} = w_t - \alpha v_{t+1}$$

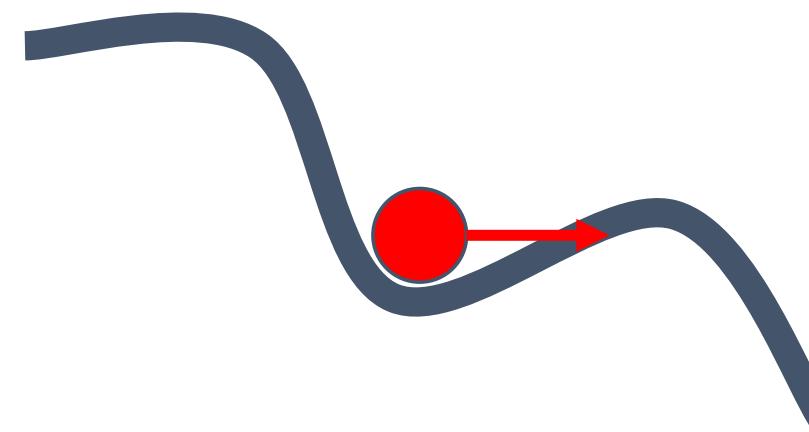
```
v = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    v = rho * v + dw
    w -= learning_rate * v
```

You may see SGD+Momentum formulated different ways,  
but they are equivalent - give same sequence of  $w$

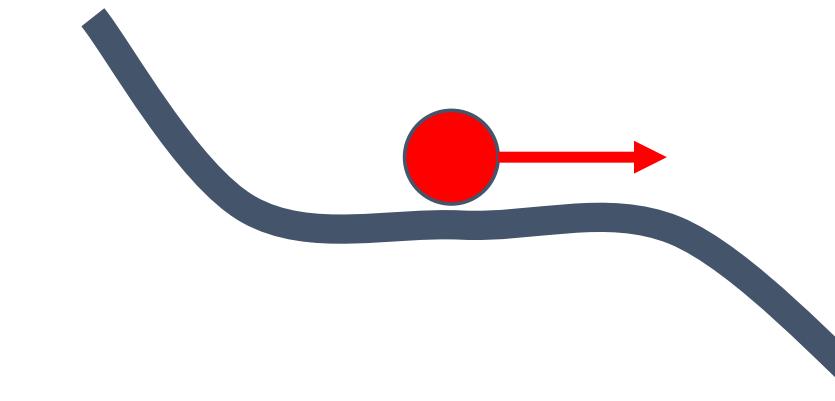


# SGD + Momentum

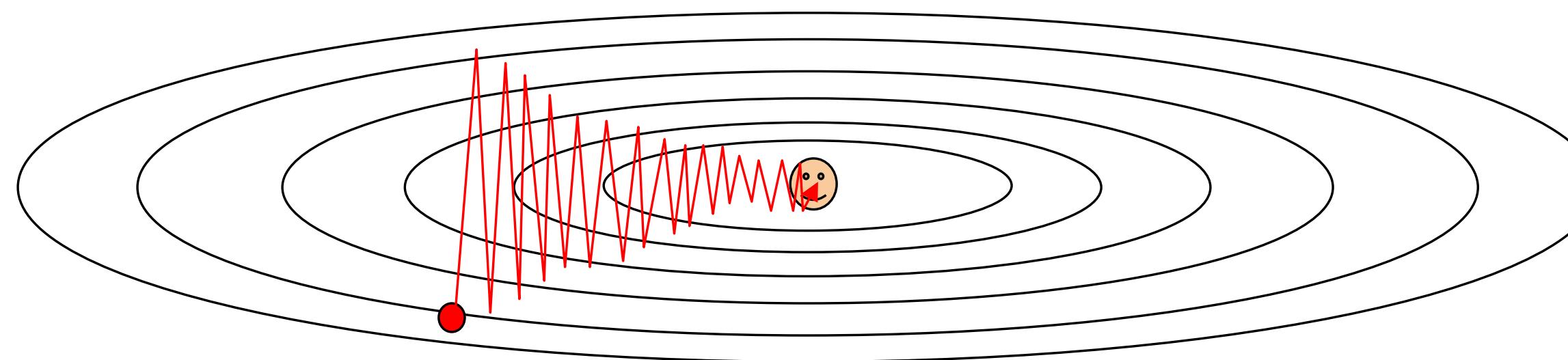
Local Minima



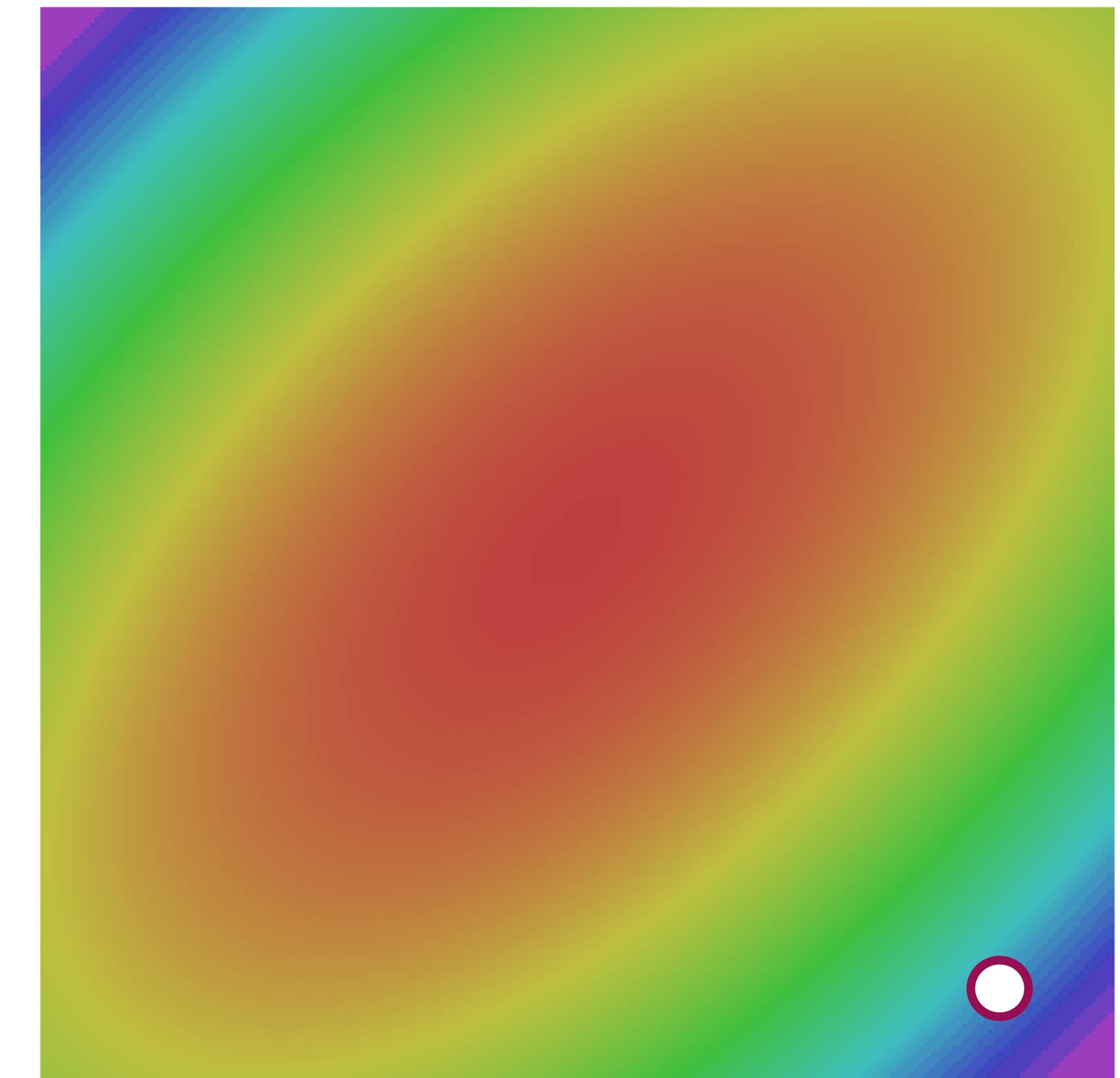
Saddle Points



Poor Conditioning



Gradient Noise

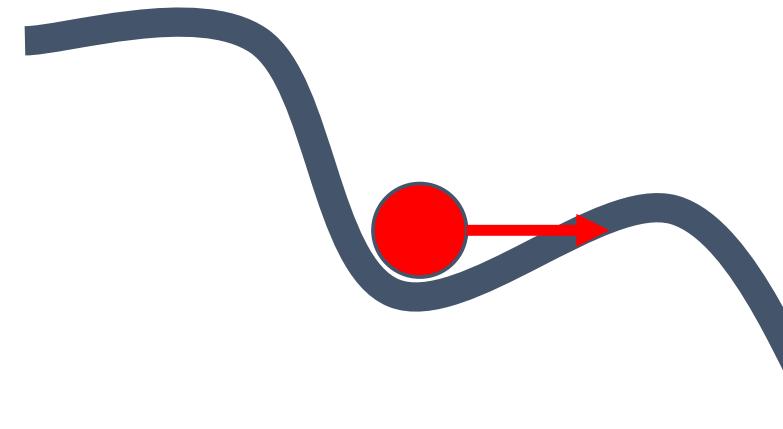


SGD

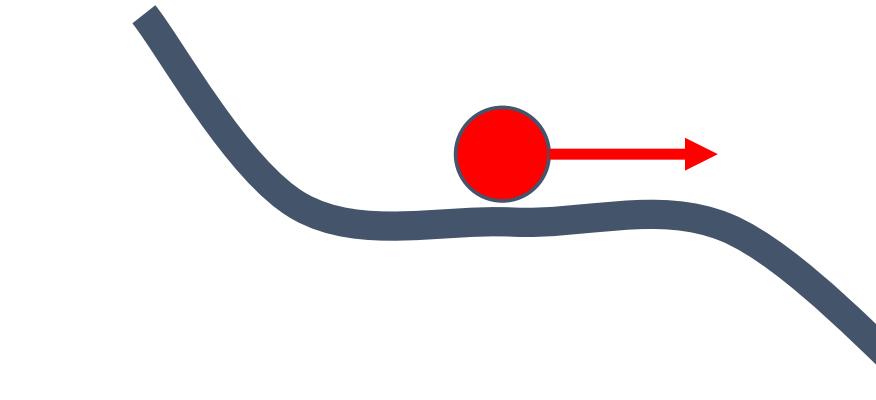
SGD+Momentum

# SGD + Momentum

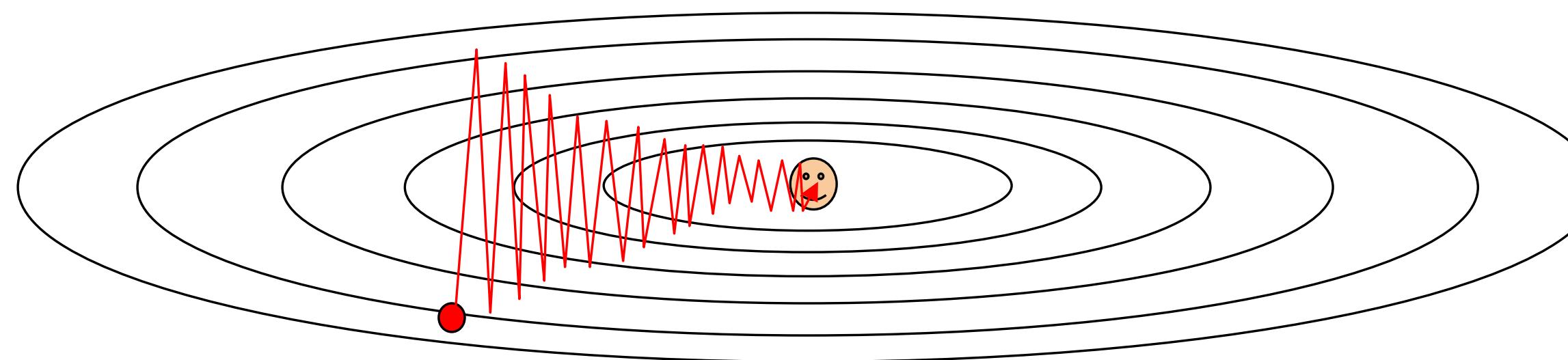
Local Minima



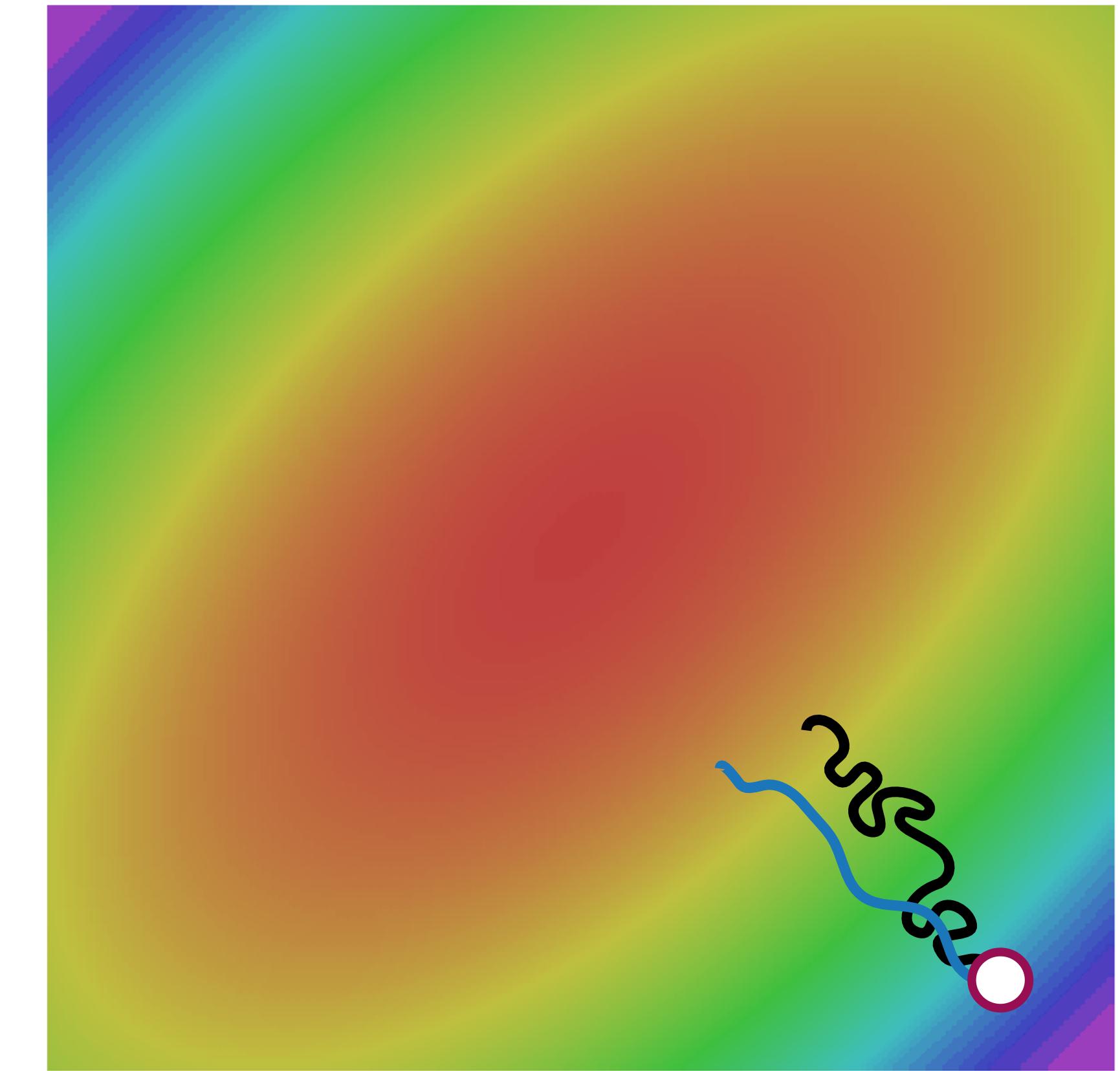
Saddle Points



Poor Conditioning



Gradient Noise

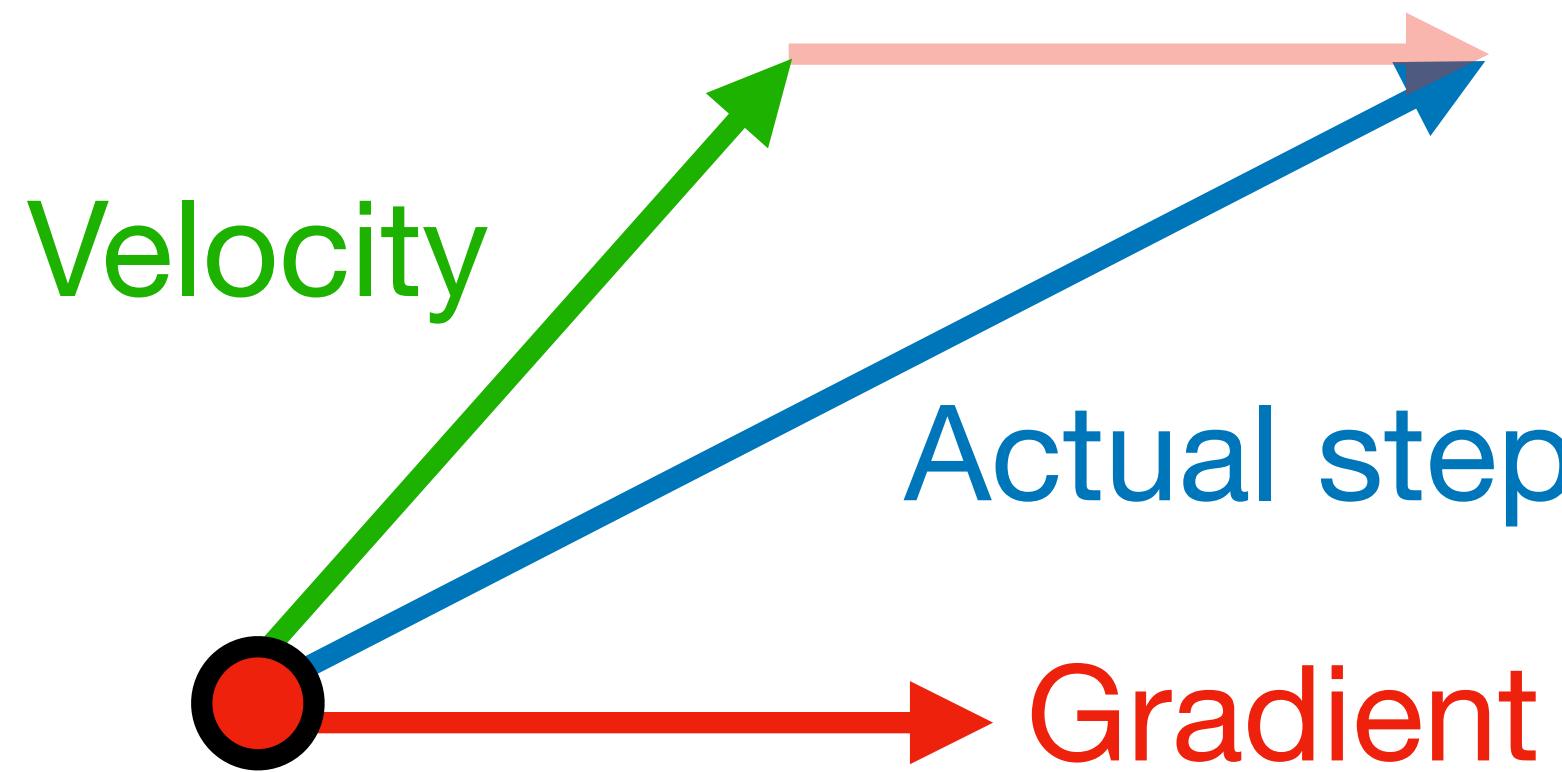


SGD

SGD+Momentum

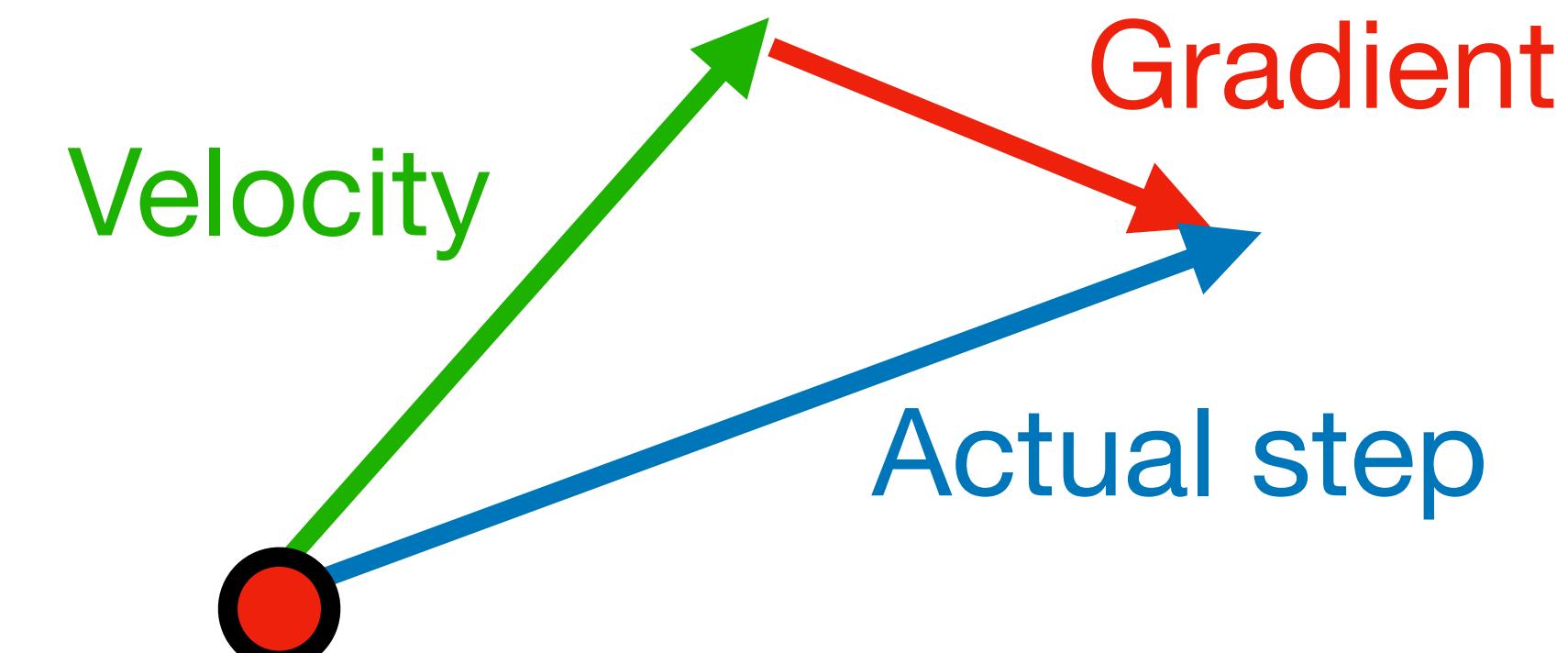
# SGD + Momentum

Momentum update:



Combine gradient at current point with velocity to get step used to update weights

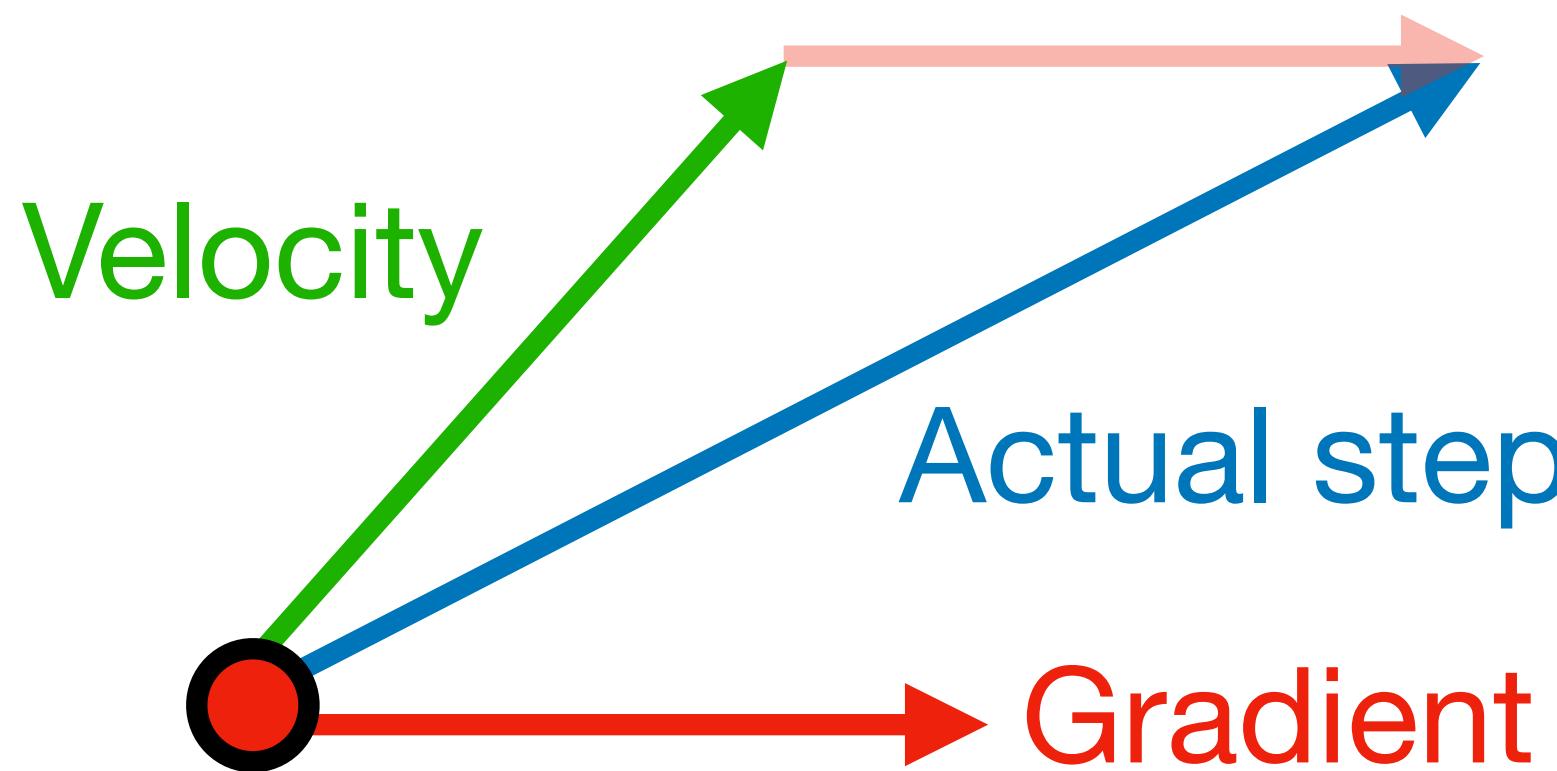
Nesterov Momentum



“Look ahead” to the point where updating using velocity would take us; compute gradient there and mix it with velocity to get actual update direction

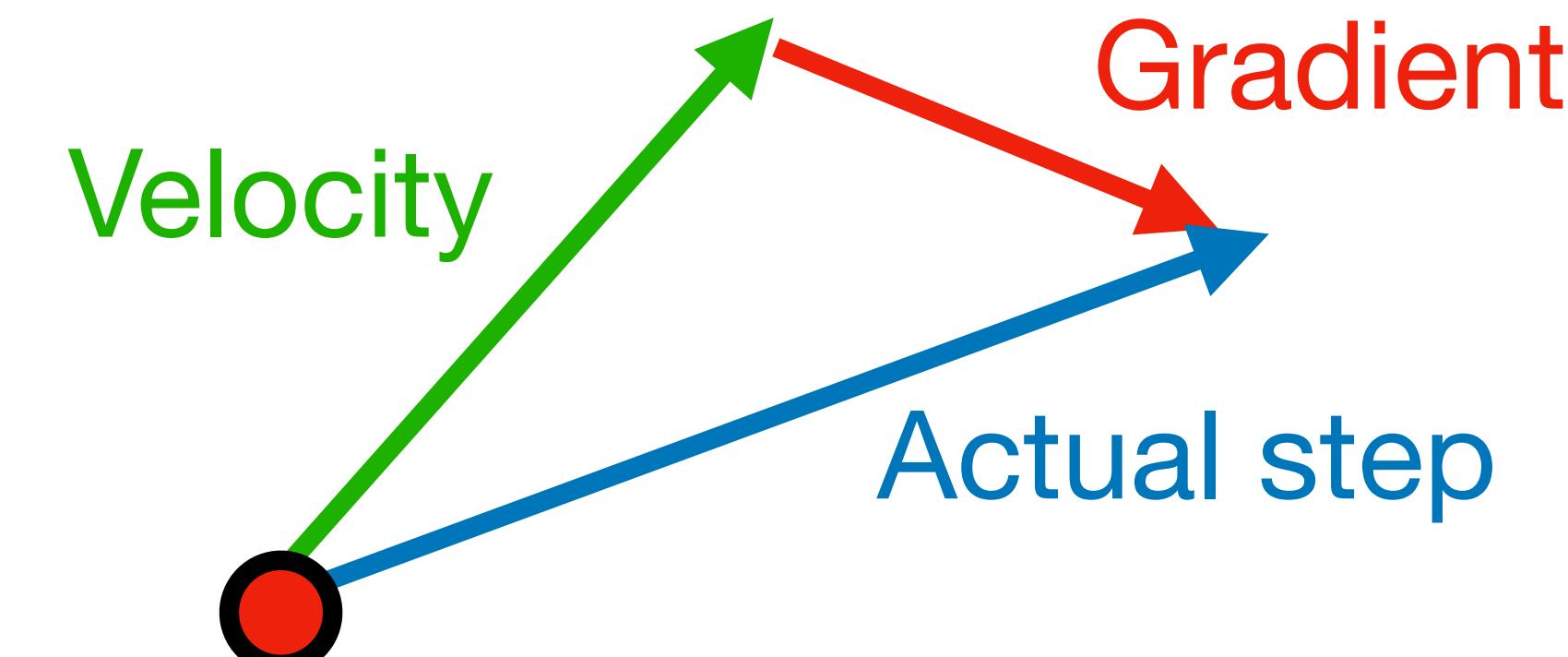
# Nesterov Momentum

Momentum update:



Combine gradient at current point with velocity to get step used to update weights

Nesterov Momentum



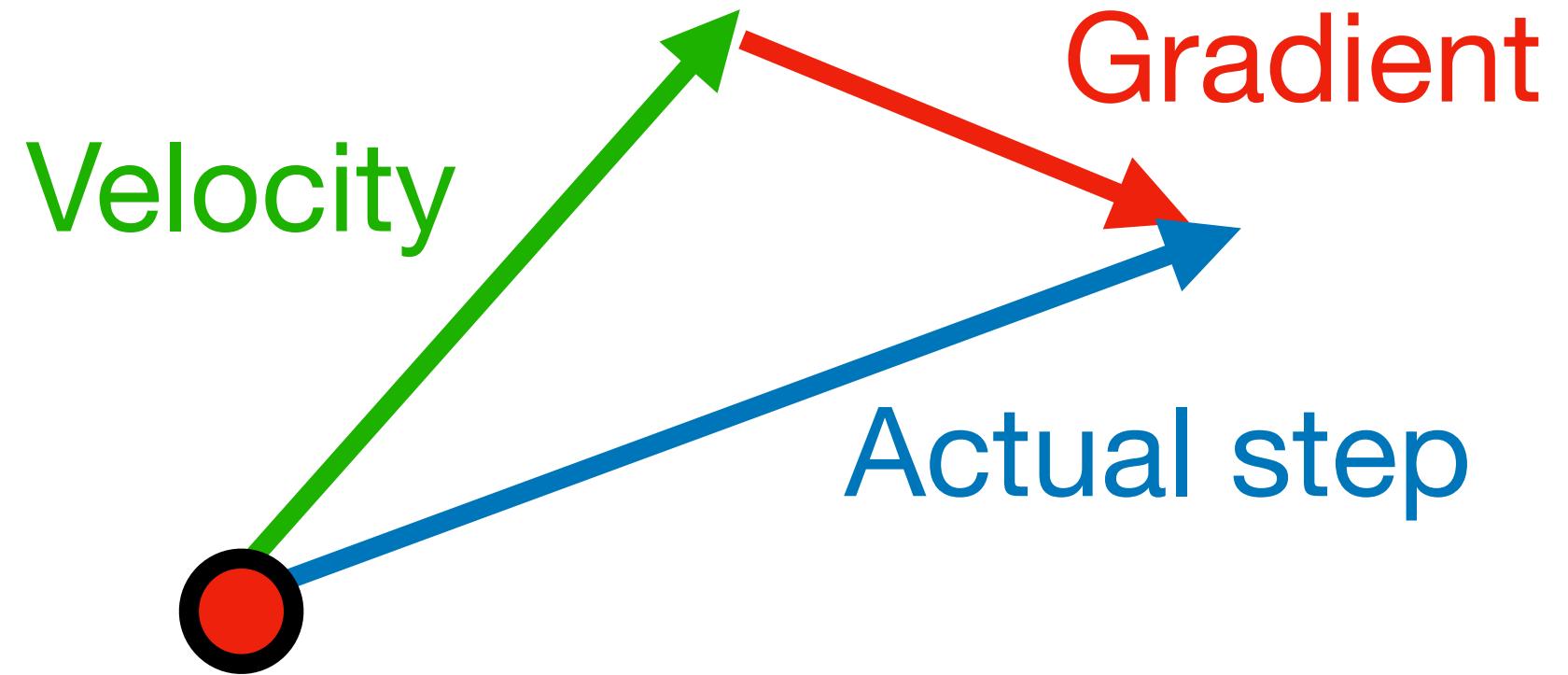
“Look ahead” to the point where updating using velocity would take us; compute gradient there and mix it with velocity to get actual update direction

# Nesterov Momentum

$$v_{t+1} = \rho v_t - \alpha \nabla L(w_t + \rho v_t)$$

$$w_{t+1} = w_t + v_{t+1}$$

Annoying, usually we want to update in terms of  $w_t, \nabla L(w_t)$



“Look ahead” to the point where updating using velocity would take us; compute gradient there and mix it with velocity to get actual update direction

# Nesterov Momentum

---

$$v_{t+1} = \rho v_t - \alpha \nabla L(w_t + \rho v_t)$$

$$w_{t+1} = w_t + v_{t+1}$$

Change of variables  
and rearrange:

$$\tilde{w}_t = w_t + \rho v_t$$

$$v_{t+1} = \rho v_t - \alpha \nabla L(\tilde{w}_t)$$

$$\tilde{w}_{t+1} = \tilde{w}_t - \rho v_t + (1 + \rho)v_{t+1}$$

$$= \tilde{w}_t + v_{t+1} + \rho(v_{t+1} - v_t)$$

Annoying, usually we want to update in terms of  $w_t, \nabla L(w_t)$

```
v = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    old_v = v
    v = rho * v - learning_rate * dw
    w -= rho * old_v - (1 + rho) * v
```



# AdaGrad

```
grad_squared = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    grad_squared += dw * dw
    w -= learning_rate * dw / (grad_squared.sqrt() + 1e-7)
```

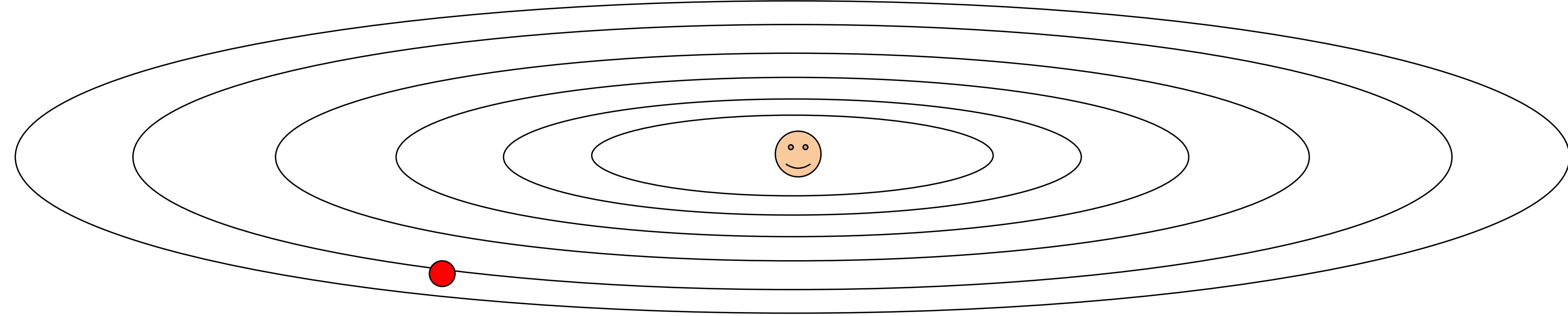
Added element-wise scaling of the gradient based on the historical sum of squares in each dimension

“Per-parameter learning rates” or “adaptive learning rates”



# AdaGrad

```
grad_squared = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    grad_squared += dw * dw
    w -= learning_rate * dw / (grad_squared.sqrt() + 1e-7)
```



**Q: What happens with AdaGrad?**

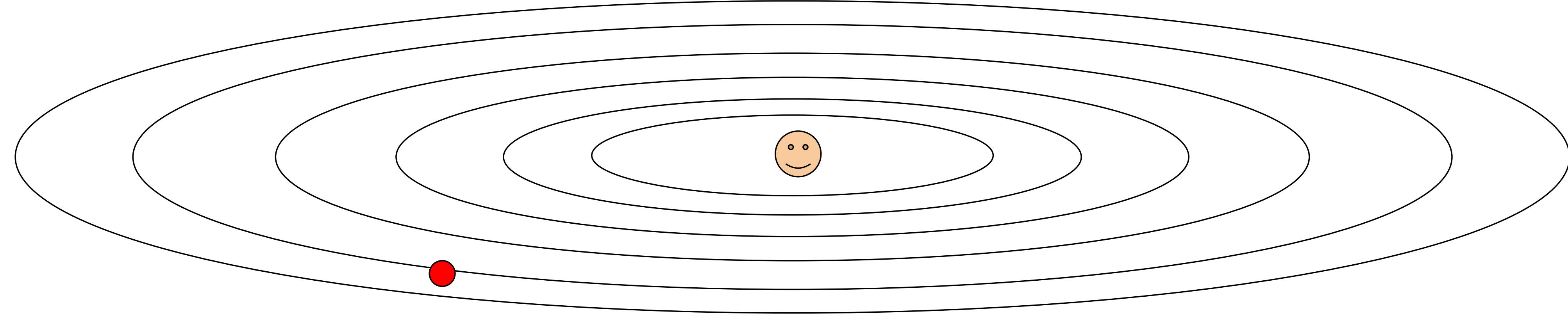
Progress along “steep” directions is damped;  
progress along “flat” directions is accelerated



# AdaGrad

```
grad_squared = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    grad_squared += dw * dw
    w -= learning_rate * dw / (grad_squared.sqrt() + 1e-7)
```

Problem: AdaGrad will slow over many iterations



**Q: What happens with AdaGrad?**

Progress along “steep” directions is damped;  
progress along “flat” directions is accelerated



# RMSProp: “Leaky AdaGrad”

```
grad_squared = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    grad_squared += dw * dw
    w -= learning_rate * dw / (grad_squared.sqrt() + 1e-7)
```

AdaGrad



```
grad_squared = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    grad_squared = decay_rate * grad_squared + (1 - decay_rate) * dw * dw
    w -= learning_rate * dw / (grad_squared.sqrt() + 1e-7)
```

RMSProp

# Adam (almost): RMSProp + Momentum

```
moment1 = 0
moment2 = 0
for t in range(1, num_steps + 1): # Start at t = 1
    dw = compute_gradient(w)
    moment1 = beta1 * moment1 + (1 - beta1) * dw
    moment2 = beta2 * moment2 + (1 - beta2) * dw * dw
    w -= learning_rate * moment1 / (moment2.sqrt() + 1e-7)
```



# Adam (almost): RMSProp + Momentum

```
moment1 = 0
moment2 = 0
for t in range(1, num_steps + 1): # Start at t = 1
    dw = compute_gradient(w)
    moment1 = beta1 * moment1 + (1 - beta1) * dw
    moment2 = beta2 * moment2 + (1 - beta2) * dw * dw
    w -= learning_rate * moment1 / (moment2.sqrt() + 1e-7)
```

Adam

Momentum

```
v = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    v = rho * v + dw
    w -= learning_rate * v
```

SGD+Momentum



# Adam (almost): RMSProp + Momentum

```
moment1 = 0
moment2 = 0
for t in range(1, num_steps + 1): # Start at t = 1
    dw = compute_gradient(w)
    moment1 = beta1 * moment1 + (1 - beta1) * dw
    moment2 = beta2 * moment2 + (1 - beta2) * dw * dw
    w -= learning_rate * moment1 / (moment2.sqrt() + 1e-7)
```

Adam

Momentum

AdaGrad / RMSProp

```
grad_squared = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    grad_squared = decay_rate * grad_squared + (1 - decay_rate) * dw * dw
    w -= learning_rate * dw / (grad_squared.sqrt() + 1e-7)
```

RMSProp



# Adam (almost): RMSProp + Momentum

```
moment1 = 0
moment2 = 0
for t in range(1, num_steps + 1): # Start at t = 1
    dw = compute_gradient(w)
    moment1 = beta1 * moment1 + (1 - beta1) * dw
    moment2 = beta2 * moment2 + (1 - beta2) * dw * dw
    w -= learning_rate * moment1 / (moment2.sqrt() + 1e-7)
```

Adam

Momentum

AdaGrad / RMSProp

Bias correction

Q: What happens at t=1?  
(Assume beta2 = 0.999)



# Adam (almost): RMSProp + Momentum

```
moment1 = 0
moment2 = 0

for t in range(1, num_steps + 1): # Start at t = 1
    dw = compute_gradient(w)
    moment1 = beta1 * moment1 + (1 - beta1) * dw
    moment2 = beta2 * moment2 + (1 - beta2) * dw * dw
    moment1_unbias = moment1 / (1 - beta1 ** t)
    moment2_unbias = moment2 / (1 - beta2 ** t)
    w -= learning_rate * moment1_unbias / (moment2_unbias.sqrt() + 1e-7)
```

Momentum

AdaGrad / RMSProp

Bias correction

**Bias correction** for the fact that first  
and second moment estimates start at  
zero

Adam with  $\beta_1 = 0.9$ ,  
 $\beta_2 = 0.999$ , and  $\text{learning\_rate} = 1e-3, 5e-4, 1e-4$   
is a great starting point for many models!



# Adam: Very common in Practice!

for input to the CNN; each colored pixel in the image yields a 7D one-hot vector. Following common practice, the network is trained end-to-end using stochastic gradient descent with the Adam optimizer [22]. We anneal the learning rate to 0 using a half cosine schedule without restarts [28].

Bakhtin, van der Maaten, Johnson, Gustafson, and Girshick, NeurIPS 2019

We train all models using Adam [23] with learning rate  $10^{-4}$  and batch size 32 for 1 million iterations; training takes about 3 days on a single Tesla P100. For each mini-batch we first update  $f$ , then update  $D_{img}$  and  $D_{obj}$ .

Johnson, Gupta, and Fei-Fei, CVPR 2018

ganized into three residual blocks. We train for 25 epochs using Adam [27] with learning rate  $10^{-4}$  and 32 images per batch on 8 Tesla V100 GPUs. We set the cubify thresh-

Gkioxari, Malik, and Johnson, ICCV 2019

sampled with each bit drawn uniformly at random. For gradient descent, we use Adam [29] with a learning rate of  $10^{-3}$  and default hyperparameters. All models are trained with batch size 12. Models are trained for 200 epochs, or 400 epochs if being trained on multiple noise layers.

Zhu, Kaplan, Johnson, and Fei-Fei, ECCV 2018

16 dimensional vectors. We iteratively train the Generator and Discriminator with a batch size of 64 for 200 epochs using Adam [22] with an initial learning rate of 0.001.

Gupta, Johnson, et al, CVPR 2018

**Adam with beta1 = 0.9, beta2 = 0.999, and learning\_rate = 1e-3, 5e-4, 1e-4 is a great starting point for many models!**



# Optimization Algorithm Comparison

Algorithm	Tracks first moments (Momentum)	Tracks second moments (Adaptive learning rates)	Leaky second moments	Bias correction for moment estimates
SGD	✗	✗	✗	✗
SGD+Momentum	✓	✗	✗	✗
Nesterov	✓	✗	✗	✗
AdaGrad	✗	✓	✗	✗
RMSProp	✗	✓	✓	✗
Adam	✓	✓	✓	✓



# L2 Regularization vs Weight Decay

## Optimization Algorithm

$$L(w) = L_{data}(w) + L_{reg}(w)$$

$$g_t = \nabla L(w_t)$$

$$s_t = \text{optimizer}(g_t)$$

$$w_{t+1} = w_t - \alpha s_t$$

L2 Regularization and Weight Decay are equivalent for SGD, SGD+Momentum so people often use the terms interchangeably!

But they are not the same for adaptive methods (AdaGrad, RMSProp, Adam, etc)

## L2 Regularization

$$L(w) = L_{data}(w) + \lambda |w|^2$$

$$g_t = \nabla L(w_t) = \nabla L_{data}(w_t) + 2\lambda w_t$$

$$s_t = \text{optimizer}(g_t)$$

$$w_{t+1} = w_t - \alpha s_t$$

## Optimization Algorithm

$$L(w) = L_{data}(w)$$

$$g_t = \nabla L_{data}(w_t)$$

$$s_t = \text{optimizer}(g_t) + 2\lambda w_t$$

$$w_{t+1} = w_t - \alpha s_t$$



# AdamW: Decouple Weight Decay

**Algorithm 2** Adam with L<sub>2</sub> regularization and Adam with decoupled weight decay (AdamW)

```

1: given  $\alpha = 0.001, \beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}, \lambda \in \mathbb{R}$ 
2: initialize time step  $t \leftarrow 0$ , parameter vector  $\theta_{t=0} \in \mathbb{R}^n$ , first moment vector  $m_{t=0} \leftarrow \theta$ , second moment
   vector  $v_{t=0} \leftarrow \theta$ , schedule multiplier  $\eta_{t=0} \in \mathbb{R}$ 
3: repeat
4:    $t \leftarrow t + 1$ 
5:    $\nabla f_t(\theta_{t-1}) \leftarrow \text{SelectBatch}(\theta_{t-1})$                                  $\triangleright$  select batch and return the corresponding gradient
6:    $\mathbf{g}_t \leftarrow \nabla f_t(\theta_{t-1}) + \lambda \theta_{t-1}$ 
7:    $\mathbf{m}_t \leftarrow \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t$                                  $\triangleright$  here and below all operations are element-wise
8:    $\mathbf{v}_t \leftarrow \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t^2$ 
9:    $\hat{\mathbf{m}}_t \leftarrow \mathbf{m}_t / (1 - \beta_1^t)$                                           $\triangleright \beta_1$  is taken to the power of  $t$ 
10:   $\hat{\mathbf{v}}_t \leftarrow \mathbf{v}_t / (1 - \beta_2^t)$                                           $\triangleright \beta_2$  is taken to the power of  $t$ 
11:   $\eta_t \leftarrow \text{SetScheduleMultiplier}(t)$                                  $\triangleright$  can be fixed, decay, or also be used for warm restarts
12:   $\theta_t \leftarrow \theta_{t-1} - \eta_t \left( \alpha \hat{\mathbf{m}}_t / (\sqrt{\hat{\mathbf{v}}_t} + \epsilon) + \lambda \theta_{t-1} \right)$ 
13: until stopping criterion is met
14: return optimized parameters  $\theta_t$ 
```



# AdamW: Decouple Weight Decay

**Algorithm 2** Adam with L<sub>2</sub> regularization and Adam with decoupled weight decay (AdamW)

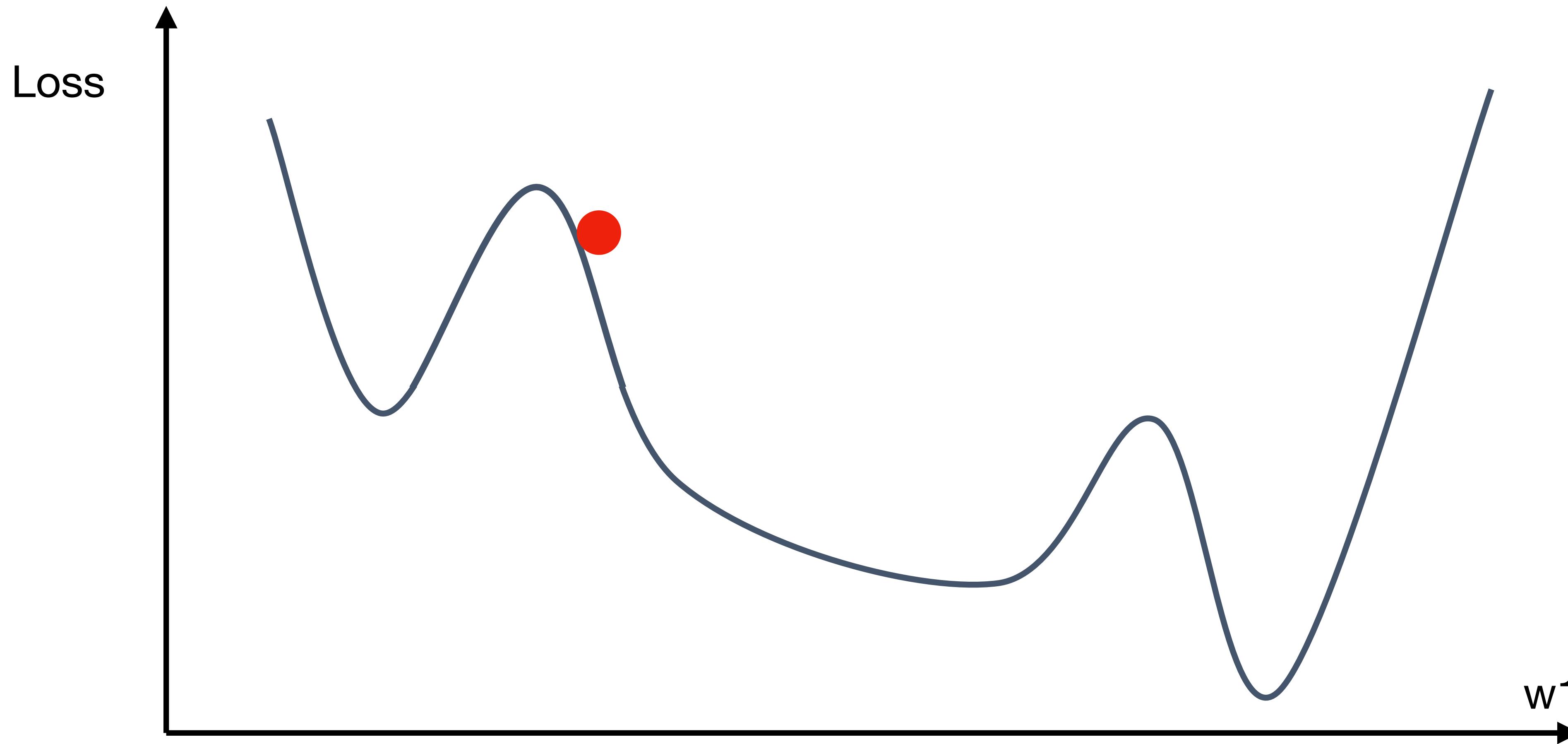
- 1: **given**  $\alpha = 0.001, \beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}, \lambda \in \mathbb{R}$
- 2: **initialize** time step  $t \leftarrow 0$ , parameter vector  $\theta_{t=0} \in \mathbb{R}^n$ , first moment vector  $m_{t=0} \leftarrow \theta$ , second moment vector  $v_{t=0} \leftarrow \theta$ , schedule multiplier  $\eta_{t=0} \in \mathbb{R}$

**AdamW should probably be your “default” optimizer for new problems**

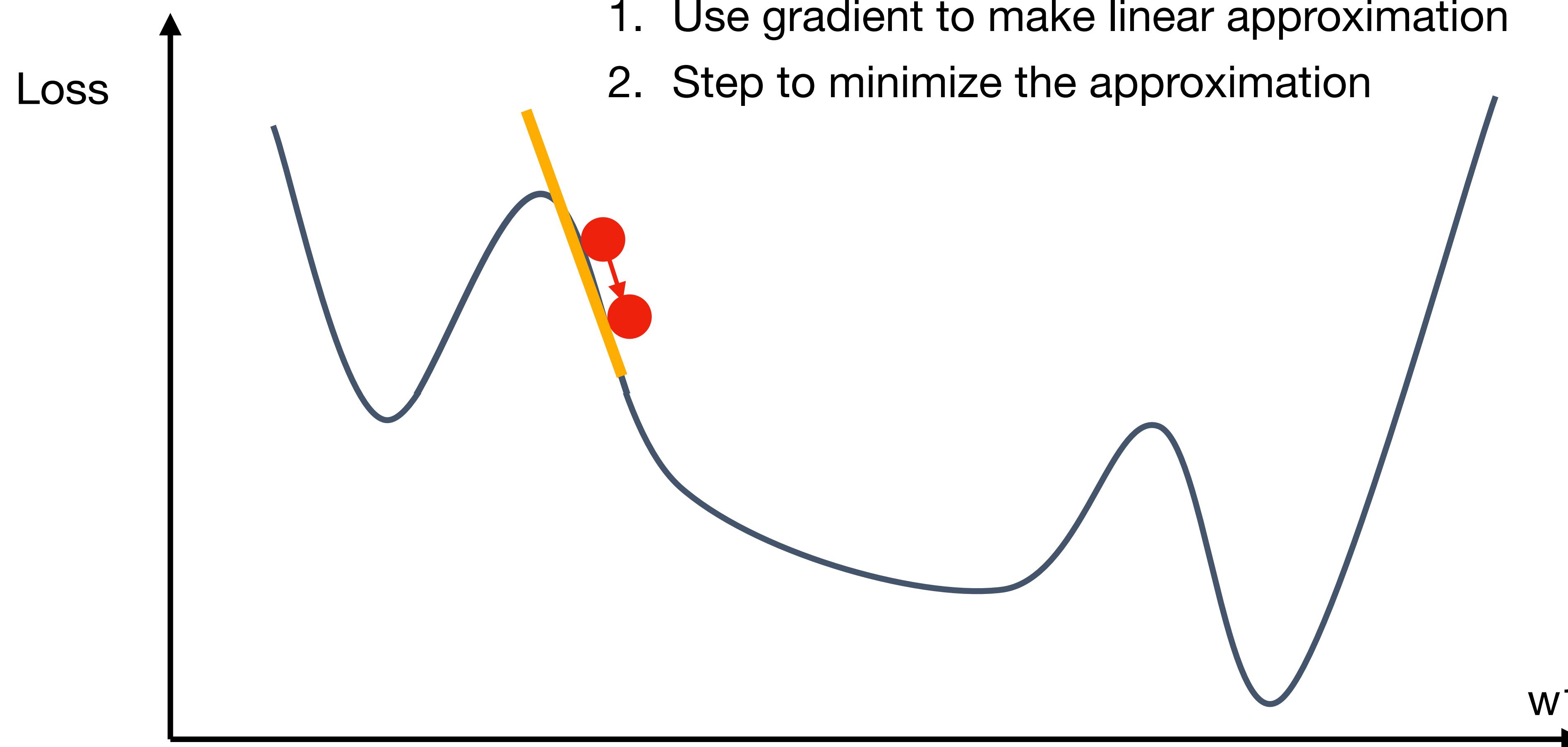
- 12:  $\theta_t \leftarrow \theta_{t-1} - \eta_t \left( \alpha \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon) + \lambda \theta_{t-1} \right)$
- 13: **until** stopping criterion is met
- 14: **return** optimized parameters  $\theta_t$



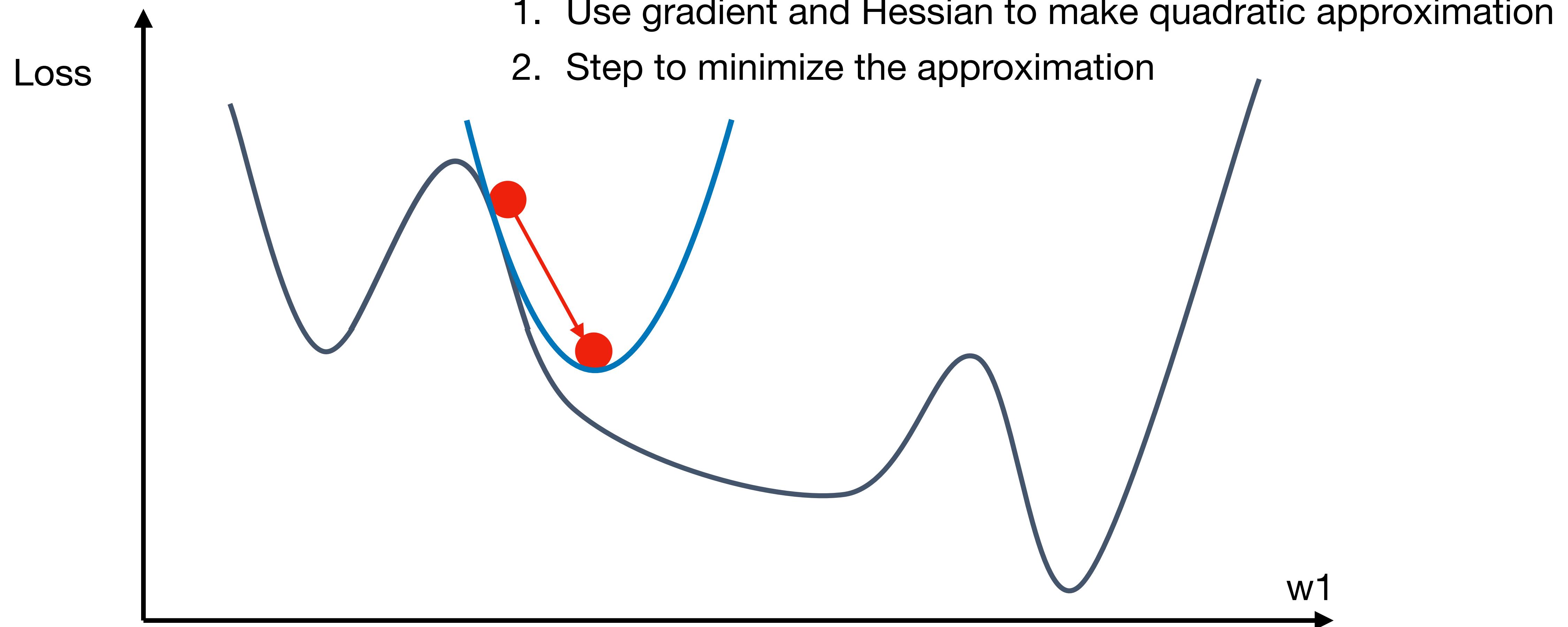
# So far: First-order Optimization



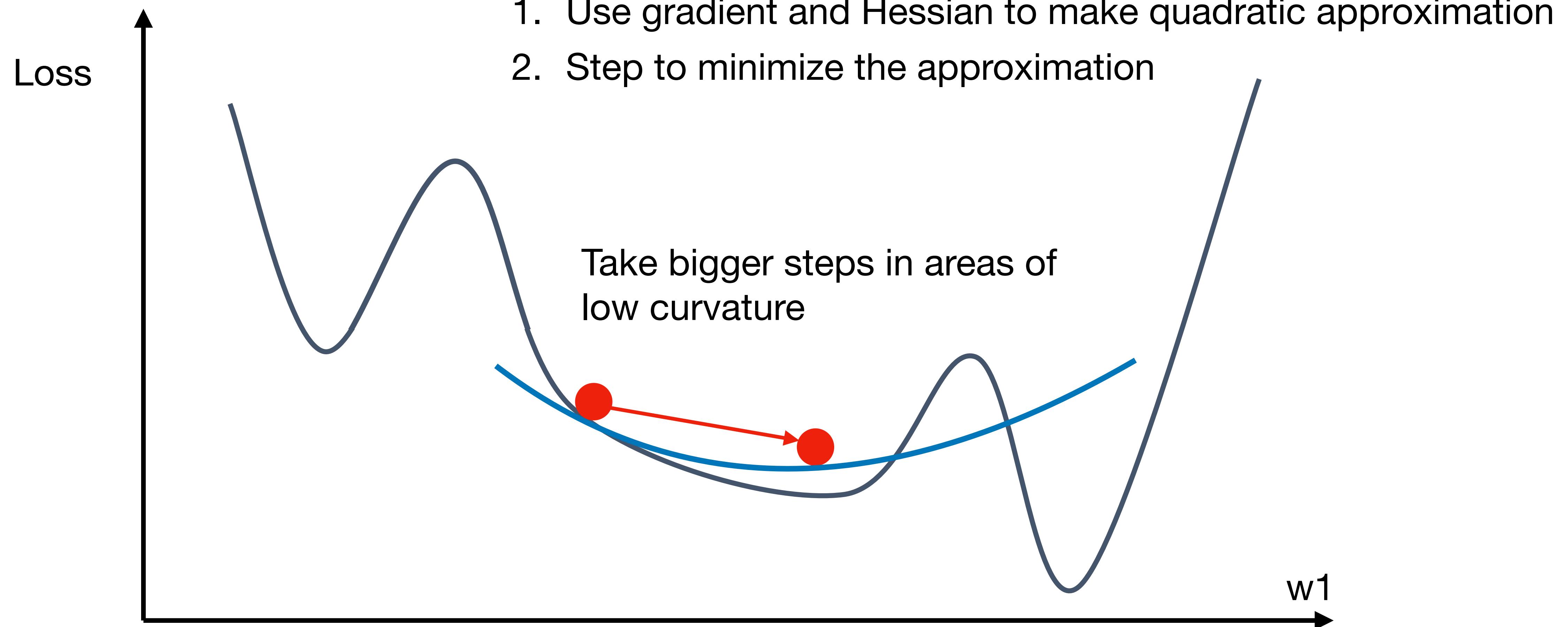
# So far: First-order Optimization



# Second-order Optimization



# Second-order Optimization



# Second-order Optimization

---

Second-order Taylor Expansion:

$$L(w) \approx L(w_0) + (w - w_0)^T \nabla_w L(w_0) + \frac{1}{2}(w - w_0)^T H_w L(w_0)(w - w_0)$$

Solving for the critical point we obtain the Newton parameter update:

$$w^* = w_0 - H_w L(w_0)^{-1} \nabla_w L(w_0)$$

# Second-order Optimization

---

Second-order Taylor Expansion:

$$L(w) \approx L(w_0) + (w - w_0)^T \nabla_w L(w_0) + \frac{1}{2}(w - w_0)^T H_w L(w_0)(w - w_0)$$

Solving for the critical point we obtain the Newton parameter update:

$$w^* = w_0 - H_w L(w_0)^{-1} \nabla_w L(w_0)$$

**Q: Why is this impractical?**

Hessian has  $O(N^2)$  elements

Inverting takes  $O(N^3)$

$N = (\text{Tens or Hundreds of Millions})$



# Second-order Optimization

---

$$w^* = w_0 - \mathbf{H}_w L(w_0)^{-1} \nabla_w L(w_0)$$

- Quasi-Newton methods (BFGS most popular): *instead of inverting the Hessian ( $O(n^3)$ ), approximate inverse Hessian with rank 1 updates over time ( $O(n^2)$  each).*
- **L-BFGS** (Limited memory BFGS): *Does not form/store the full inverse Hessian*

# Second-order Optimization: L-BFGS

---

- **Usually works very well in full batch, deterministic mode** i.e. if you have a single, deterministic  $f(x)$  then L-BFGS will probably work very nicely.
- **Does not transfer very well to mini-batch setting.** Gives bad results. Adapting second-order methods to large-scale, stochastic setting is an active area of research.



# In practice:

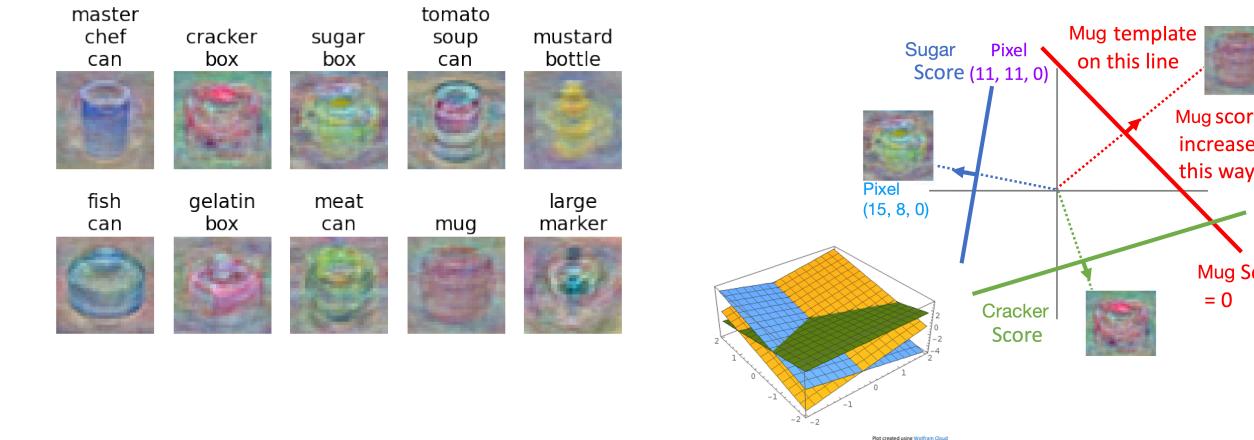
---

- **Adam** is a good default choice in many cases  
**SGD+Momentum** can outperform Adam but may require more tuning.
- If you can afford to do full batch updates then try out **L-BFGS** (and don't forget to disable all sources of noise)

# Summary

- Use **Linear Models** for image classification problems.
- Use **Loss Functions** to express preferences over different choices of weights.
- Use **Regularization** to prevent overfitting to training data.
- Use **Stochastic Gradient Descent** to minimize our loss functions and train the model.

$$s = f(x; W) = Wx$$



$$L_i = -\log\left(\frac{\exp^{s_{y_i}}}{\sum_j \exp^{s_j}}\right)$$

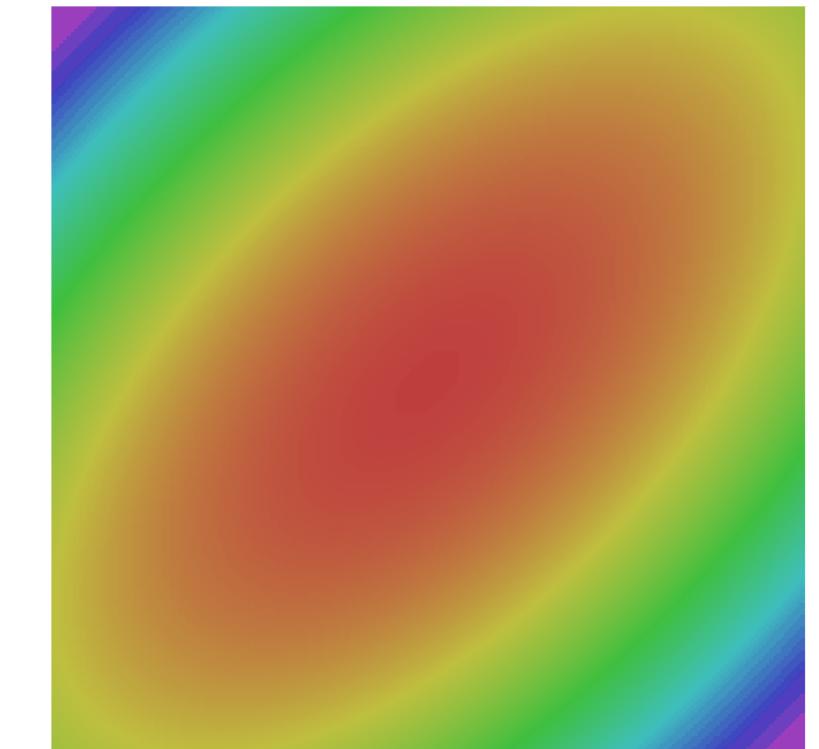
**Softmax**

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

**SVM**

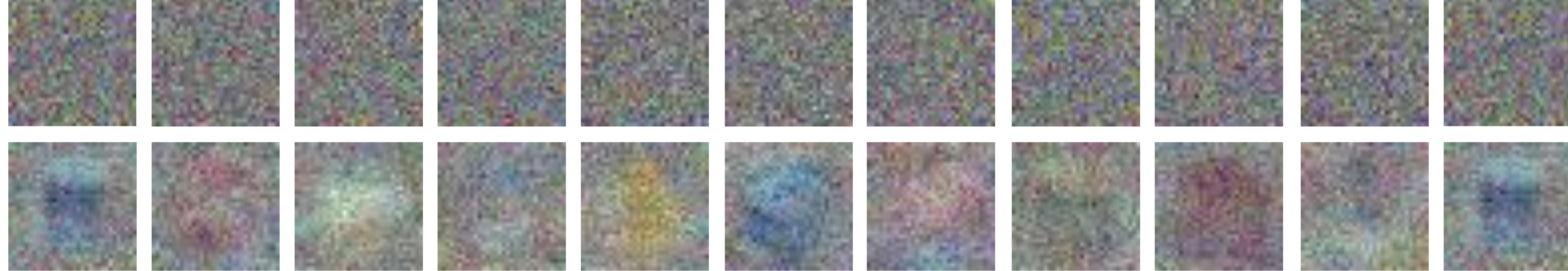
$$L = \frac{1}{N} \sum_{i=1}^N L_i + R(W)$$

```
v = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    v = rho * v + dw
    w -= learning_rate * v
```



# Next time: Neural Networks

**DR**



# DeepRob

**Lecture 4**  
**Regularization + Optimization**  
**University of Michigan and University of Minnesota**

